

# Non-Destructive Prediction of Fruit Ripeness and Firmness Using Hyperspectral Imaging and Lightweight Machine Learning Models

Phongsakon Mark Konrad<sup>a</sup>, Casper Kunstmann-Olsen<sup>b</sup>, Jacek Fiutowski<sup>b</sup>, Serkan Ayvaz<sup>a,\*</sup>

<sup>a</sup>*The Maersk Mc-Kinney Moller Institute, SDU Centre for Industrial Software, University of Southern Denmark, Sønderborg, 6400, Denmark*

<sup>b</sup>*Mads Clausen Institute, SDU NanoSyd, University of Southern Denmark, Sønderborg, 6400, Denmark*

---

## Abstract

Post-harvest fruit quality assessment is essential for reducing food waste, yet reliable non-destructive methods typically depend on expensive hyperspectral cameras and computationally intensive deep learning models. These systems typically require GPU resources, large-scale training data, and domain expertise, limiting their feasibility for many real-world agricultural settings. This study systematically evaluates 20 classical machine learning algorithms on hyperspectral imaging data for simultaneous ripeness classification and firmness prediction across five fruit species, using cross-validated experimental design with Bayesian hyperparameter optimization. Data preprocessing strategy, particularly class balancing and spectral transformations, contributes as much to prediction accuracy as algorithm choice. Our results show that tree-based machine learning models can outperform state-of-the-art deep

---

\*Corresponding author

*Email address:* seay@mimi.sdu.dk (Serkan Ayvaz)

learning models reported in Fruit-HSNet. Moreover, the findings indicate that only three visible-range wavelengths are needed to recover over 94% of full-spectrum accuracy, demonstrating that low-cost multispectral sensors combined with lightweight machine learning models can serve as practical alternatives to expensive hyperspectral cameras and complex deep learning approaches for practical fruit quality sorting.

*Keywords:* Hyperspectral imaging, Machine learning, Multi-attribute Fruit quality assessment, Non-destructive testing, Precision agriculture

---

## 1. Introduction

Commercial fruit quality assessment relies on destructive sampling of small fractions of harvests or subjective visual inspection (Lu et al., 2020). A typical processing facility handles tens of thousands of items daily but can destructively test perhaps 100 samples. The remainder is graded by surface appearance, a weak proxy for internal quality that leads to misclassification, post-harvest losses, and inefficient supply-chain decisions.

Hyperspectral imaging (HSI) addresses this gap through non-destructive, contact-free measurement of wavelength-dependent reflectance across the visible and near-infrared spectrum. Spectral signatures encode biochemical and structural changes associated with ripeness, firmness, and spoilage (Yang et al., 2025). Laboratory studies report accuracy exceeding 84% for strawberry maturity assessment (Su et al., 2021), yet commercial adoption remains limited.

Deep learning architectures achieve strong laboratory performance: 3D-CNNs reach 84% accuracy for strawberry classification (Su et al., 2021),

17 transformer models reach 99.4% for contamination detection (Guo et al.,  
18 2024), and real-time systems achieve 98.6% accuracy (Gao et al., 2020). How-  
19 ever, these approaches require GPU acceleration, large training datasets, and  
20 specialist expertise that most agricultural operations cannot justify. Classical  
21 machine learning offers a different set of tradeoffs. These methods train effi-  
22 ciently on limited data, run on consumer-grade CPU hardware, and produce  
23 interpretable feature importance scores amenable to validation against plant  
24 physiology. Whether classical methods can achieve competitive accuracy on  
25 hyperspectral fruit quality benchmarks remains uncharacterized.

26 Systematic benchmarking is lacking. Published comparisons typically fix  
27 a single preprocessing pipeline without ablation, leaving open the question  
28 of how much accuracy variation stems from data engineering rather than  
29 algorithm selection.

30 This study benchmarks 20 classical and gradient-boosted machine learn-  
31 ing algorithms on paired ripeness and firmness prediction using the DeepHS  
32 Fruit dataset (Varga et al., 2021). A six-phase experimental design spans  
33 200 model-configuration pairs, 4,000 Bayesian optimization trials, 10-fold  
34 cross-validation, ensemble learning, and explainable AI analysis. The central  
35 question is whether preprocessing decisions contribute as much to classifi-  
36 cation accuracy as algorithm choice, and if so, what mechanisms drive that  
37 effect. By evaluating all 20 algorithms under identical, controlled conditions,  
38 this benchmark provides a transparent reference point for the field. The main  
39 contributions are:

- 40 1. Systematic ablation across 200 configurations provides the first evi-  
41 dence that preprocessing shifts accuracy as much as algorithm selection

- 42 (25.2 vs. 32.4 pp; bootstrap 95% CI on the difference contains zero).
- 43 2. Firmness prediction outperforms ripeness classification across all 20  
44 models ( $d=1.84$ ), traced to label-structure asymmetry rather than algorithm-  
45 specific effects.
  - 46 3. PCA degrades performance after engineered spectral transforms be-  
47 cause variance-maximizing compression discards low-variance features  
48 that carry discriminative information.
  - 49 4. Consensus feature ranking identifies three visible-range wavelengths  
50 (448, 540, 640 nm) that recover 94.7% of full-spectrum accuracy; ap-  
51 proximate RGB wavelengths (450, 550, 651 nm) perform comparably  
52 (96.1%), indicating that band count rather than precise band place-  
53 ment drives the reduced-band performance.

54 The remainder of this paper is organized as follows. Section 2 reviews  
55 prior work on HSI-based fruit quality assessment, preprocessing pipelines,  
56 and explainability methods. Section 3 describes the dataset, feature engi-  
57 neering, and six-phase experimental design. Section 4 presents results across  
58 all phases, including ablation, cross-validation, and XAI analysis. Section 5  
59 interprets the key findings and discusses limitations. Section 6 summarizes  
60 the contributions and outlines directions for future work.

## 61 **2. Related Work**

62 Three areas of prior work frame the contribution of this study.

### 63 *2.1. HSI and Machine Learning for Fruit Quality*

64 Hyperspectral systems acquire spectral information across hundreds of  
65 narrow wavelength bands, spanning visible (400–700 nm) and near-infrared

66 (700–2500 nm) regions. This spectral data enables detection of biochemical  
67 changes during fruit maturation: chlorophyll degradation, carotenoid accu-  
68 mulation, and cellular structure modifications (Lu et al., 2020). Spectral-  
69 biochemical relationships have been established for quality attributes: chloro-  
70 phyll absorption peaks around 680 nm indicate ripeness state, while near-  
71 infrared signatures (850–950 nm) correlate with cellular structure and firm-  
72 ness (Khodabakhshian and Emadi, 2017; Feng et al., 2023). Translating these  
73 laboratory findings into deployable systems remains difficult.

74 Traditional machine learning approaches, including Partial Least Squares  
75 Regression (PLSR), Support Vector Machines (SVM), and tree-based ensem-  
76 bles, achieved 70–85% accuracy with advantages in interpretability and com-  
77 putational efficiency (Lu et al., 2017; Khodabakhshian and Emadi, 2017).  
78 Reviews of the field identify PLSR and SVM as the most commonly ap-  
79 plied methods (Wieme et al., 2022), and targeted wavelength selection can  
80 maintain classification accuracy with as few as six bands (Nagasubramanian  
81 et al., 2018). These practical insights have received limited attention as the  
82 field shifted toward deep learning, and systematic evaluation of traditional  
83 methods under standardized conditions remains sparse.

84 Despite these accuracy gains, deep learning models require GPU acceler-  
85 ation, specialized expertise, and extensive training datasets unavailable for  
86 diverse fruit varieties or regional conditions. Computational requirements  
87 exceed what most agricultural facilities can justify. Hybrid approaches com-  
88 bining traditional feature engineering with deep learning (Liu et al., 2024;  
89 Olisah et al., 2024) often compound complexity without addressing deploy-  
90 ment barriers.

91 *2.2. Multi-Attribute Prediction and Explainability*

92 Ripeness and firmness are biologically coupled through ethylene produc-  
93 tion, cell wall breakdown, and chlorophyll degradation (Feng et al., 2023;  
94 Varga et al., 2021), yet most research treats them as independent prediction  
95 problems. Feng et al. (2023) found advantages for joint prediction over inde-  
96 pendent models in loquat quality assessment, and multi-modal approaches  
97 combining hyperspectral and RGB data show 10–15% accuracy improve-  
98 ments (Garillos-Manliguez and Chiang, 2021). Yet single-task evaluation  
99 persists despite the fact that isolated quality attributes rarely drive com-  
100 mercial sorting decisions. This study adopts a paired single-task evaluation  
101 framework assessing ripeness and firmness jointly via aggregated metrics,  
102 while training separate models per task.

103 SHAP and LIME analyses identify biologically meaningful wavelengths:  
104 680 nm for chlorophyll absorption, 760 nm for red-edge transitions, 970 nm  
105 for water content (Ahmed et al., 2024b). These findings align with plant  
106 physiology knowledge. If only 6–15 wavelengths contain most discrimina-  
107 tive information (Nagasubramanian et al., 2018; Ahmed et al., 2024a), full-  
108 spectrum systems may be overengineered and targeted multispectral sensors  
109 could achieve comparable performance at reduced cost.

110 *2.3. Benchmarking Gaps and Deployment Challenges*

111 Most studies evaluate 2–5 algorithms under carefully selected conditions  
112 (Yang et al., 2025). Reported accuracies range from 73% to 99% for similar  
113 classification tasks (Vignati et al., 2023). The DeepHS Fruit dataset (Varga  
114 et al., 2021) provides standardized data across five fruit varieties (avocado,

115 kiwi, mango, kaki/persimmon, papaya) with ripeness and firmness annota-  
116 tions. Ben Jmaa et al. (2025) introduced Fruit-HSNet, reporting 70.73%  
117 overall accuracy on this benchmark. Systematic evaluation of traditional  
118 machine learning methods on this benchmark remains limited, with most  
119 comparative studies focusing on deep learning architectures.

120 Hyperspectral systems cost \$10,000–100,000 compared to \$500–5,000 for  
121 RGB alternatives, require specialized operators, and experience 15–30% accu-  
122 racy degradation under variable field conditions (Min et al., 2023). Process-  
123 ing speed requirements for sorting applications (under 100 ms per sample)  
124 eliminate most deep learning approaches. RGB-reconstructed hyperspec-  
125 tral images (Ahmed et al., 2024a) and edge computing solutions (Lanke and  
126 Chandak, 2025) address symptoms rather than causes.

127 No existing study systematically benchmarks a broad set of traditional  
128 ML algorithms on multi-attribute fruit quality prediction with controlled  
129 preprocessing ablation.

### 130 **3. Materials and Methods**

#### 131 *3.1. Dataset Description*

132 We used the DeepHS Fruit dataset (version 2), a publicly available hyper-  
133 spectral imaging dataset for fruit quality assessment research (Varga et al.,  
134 2021). The dataset is hosted on GitHub<sup>1</sup> and provides hyperspectral data  
135 with quality annotations for benchmarking. Version 2 corrects a potential  
136 test-set contamination issue in the original split by ensuring that balanced

---

<sup>1</sup>[https://github.com/cogsys-tuebingen/deephs\\_fruit](https://github.com/cogsys-tuebingen/deephs_fruit)

137 training variants do not alter the benchmark test partition; train, validation,  
138 and test sets are validated for zero sample-level overlap.

139 We use the dataset-provided VIS benchmark split as our fixed evalua-  
140 tion setting (unbalanced train/test: train  $n=381$ , test  $n=138$ ; total  $n=519$ ).  
141 Fruit-type distribution in this split is avocado ( $n=170$ , 32.8%), kiwi ( $n=162$ ,  
142 31.2%), mango ( $n=68$ , 13.1%), kaki/persimmon ( $n=68$ , 13.1%), and papaya  
143 ( $n=51$ , 9.8%). The VIS subset spans 47 acquisition days, providing temporal  
144 variation in fruit maturation states. For the stratified resplit condition, we  
145 keep the benchmark test set fixed but replace the training set with a fruit-  
146 balanced training set ( $n=414$ ; avocado  $n=130$ , kiwi  $n=130$ , mango  $n=56$ ,  
147 kaki/persimmon  $n=55$ , papaya  $n=43$ ), improving balance without breaking  
148 comparability.

149 The full DeepHS Fruit dataset includes multiple hyperspectral imaging  
150 systems (Specim FX10, INNO-SPEC Redeye 1.7, and Corning microHSI 410  
151 Vis-NIR). To maintain spectral consistency and comparability with existing  
152 benchmarks, this study restricts evaluation to the VIS camera (Specim FX10,  
153 398–1004 nm), and all reported results reflect this sensor and its wavelength  
154 range. The FX10 covers only the short-NIR region; true NIR water absorp-  
155 tion bands (1200–1450 nm) that are commonly used in firmness assessment  
156 are absent from these data.

157 Each sample includes quality annotations from destructive measurements.  
158 Ripeness classification follows a three-class scheme based on visual assess-  
159 ment and destructive sampling (Varga et al., 2021); inter-annotator agree-  
160 ment was not reported in the original dataset publication. Classes are: unripe  
161 ( $n=116$ , 22.4%), perfect ( $n=273$ , 52.6%), and overripe ( $n=130$ , 25.0%). Per-

162 fect ripeness samples are overrepresented at 52.6%, while unripe and overripe  
163 samples are more evenly distributed.

164 Firmness measurements (grams-force from penetrometer testing) are cat-  
165 egorized into three classes derived from cross-fruit distribution analysis: soft  
166 (0–1000 gf,  $n=220$ , 42.4%), medium (1001–2500 gf,  $n=186$ , 35.8%), and firm  
167 (2501+ gf,  $n=99$ , 19.1%). The distribution skews toward softer fruits (42.4%  
168 soft vs 19.1% firm). Samples with missing firmness measurements ( $n=14$ ,  
169 2.7%; train  $n=13$ , test  $n=1$ ) are retained and encoded as a separate “Un-  
170 known” label during training to avoid discarding data; they are excluded  
171 from all quality-metric computations reported in the tables.

172 Sample distribution shows fruit-type imbalances (avocado  $n=170$  vs pa-  
173 paya  $n=51$ ) and ripeness-class imbalances (perfect 52.6% vs unripe 22.4%,  
174 Fig. 1a). PCA analysis indicates that 95% cumulative explained variance  
175 is achieved with approximately 18–20 principal components from the 1,120-  
176 dimensional feature space (Fig. 1b). Projection onto the first two principal  
177 components (Fig. 1c) reveals that cluster structure is driven primarily by  
178 fruit type rather than quality state. This species-dominated variance struc-  
179 ture motivates per-task rather than joint classification approaches.

## 180 *3.2. Data Preprocessing and Feature Engineering*

### 181 *3.2.1. Spectral Feature Extraction*

182 The DeepHS Fruit dataset provides hyperspectral data from three imag-  
183 ing systems (Specim FX10, INNO-SPEC Redeye 1.7, Corning microHSI 410  
184 Vis-NIR). We use only the VIS camera (Specim FX10; 224 bands, 398–  
185 1004 nm) for spectral consistency. We consume the dataset-provided ex-  
186 tracted spectral feature tables (derived from background-masked hyperspec-

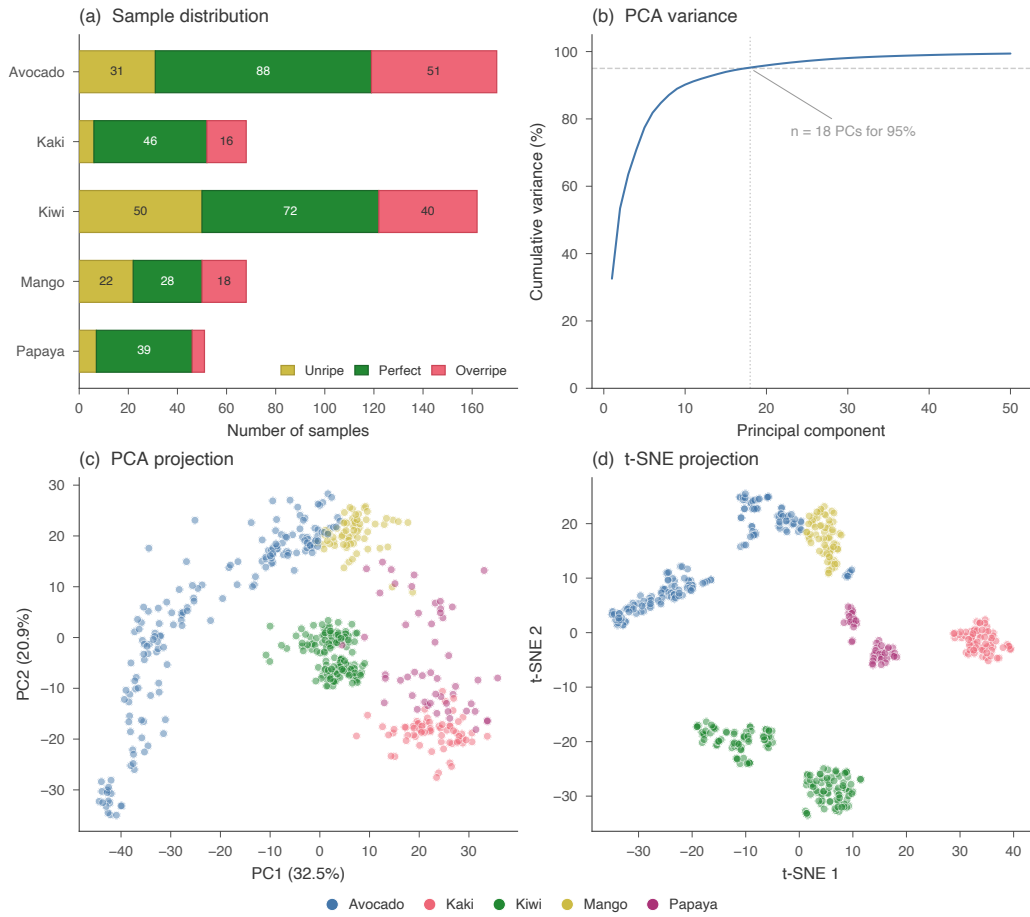


Figure 1: Dataset overview. **(a)** Sample distribution across five fruit types and three ripeness classes. **(b)** PCA cumulative variance: 18 components capture 95% of variance in the 1,120-dimensional feature space. **(c)** PCA projection and **(d)** t-SNE projection colored by fruit type show species-driven cluster structure rather than quality-state separation.

187 tral cubes to isolate fruit pixels) and do not re-run pixel-level segmentation  
188 within the benchmarking pipeline.

189 An average spectrum is calculated for each sample by averaging spectral  
190 values across all foreground (fruit) pixels for each wavelength band:

$$\bar{S}(\lambda) = \frac{1}{|\mathcal{F}|} \sum_{(i,j) \in \mathcal{F}} S_{i,j}(\lambda) \quad (1)$$

191 where  $S_{i,j}(\lambda)$  represents the spectral value at wavelength  $\lambda$  for pixel  $(i, j)$ ,  
192 and  $\mathcal{F}$  denotes the set of foreground pixels. Each hyperspectral cube reduces  
193 to a 224-dimensional spectral signature per sample.

### 194 3.2.2. Derived Spectral Features

195 The combined feature set concatenates five spectral representations from  
196 the extracted feature tables, totaling 1,120 dimensions (5 representations  $\times$   
197 224 spectral bands).

198 The average spectrum  $\bar{S}(\lambda)$  captures absolute spectral reflectance across  
199 the VIS–NIR range. The first derivative highlights regions of rapid spectral  
200 change, emphasizing transitions related to chlorophyll degradation, carotenoid  
201 accumulation, and other biochemical changes:

$$D_1(\lambda) = \frac{d\bar{S}(\lambda)}{d\lambda}. \quad (2)$$

202 Continuum removal normalizes the spectrum by fitting a convex hull and  
203 dividing the original spectrum by this hull, highlighting absorption bands  
204 such as chlorophyll absorption around 680 nm:

$$CR(\lambda) = \frac{\bar{S}(\lambda)}{C(\lambda)}, \quad (3)$$

205 where  $C(\lambda)$  is the continuum value at wavelength  $\lambda$ .

206 The standard normal variate (SNV) standardizes each mean spectrum:

$$SNV(\lambda) = \frac{\bar{S}(\lambda) - \mu}{\sigma}, \quad (4)$$

207 where  $\mu$  and  $\sigma$  are the mean and standard deviation of  $\bar{S}(\lambda)$ . Applied to  
208 per-sample mean spectra (averaged across pixels), SNV corrects for baseline  
209 offset and gain differences between samples; it does not correct within-pixel  
210 multiplicative scatter in the sense used for raw reflectance cubes.

211 Finally, the first derivative of the continuum-removed spectrum combines  
212 absorption-emphasizing and edge-enhancing properties:

$$D_1^{CR}(\lambda) = \frac{dCR(\lambda)}{d\lambda}. \quad (5)$$

213 These five representations capture complementary aspects of spectral  
214 variation.

215 Raw spectral reflectance signatures for the five fruit types exhibit char-  
216 acteristic pigment-related variations across ripeness states, particularly in  
217 chlorophyll absorption regions ( $\sim 680$  nm) and red-edge transitions ( $\sim 700$ –  
218  $750$  nm). Near-infrared features ( $700$ – $1000$  nm) show structure-related varia-  
219 tions correlated with firmness. These signatures arise from changes in cellular  
220 architecture and water content during maturation (Fig. 2).

### 221 3.2.3. Dataset Balancing and Train/Test Splits

222 The original DeepHS Fruit splits show class imbalances for certain fruit  
223 types, ripeness states, and firmness categories. We distinguish two categories  
224 of data balancing:

225 *Training set replacement.* The stratified resplit condition replaces the  
226 training partition with a fruit-balanced set ( $n=414$  vs the original  $n=381$ ),

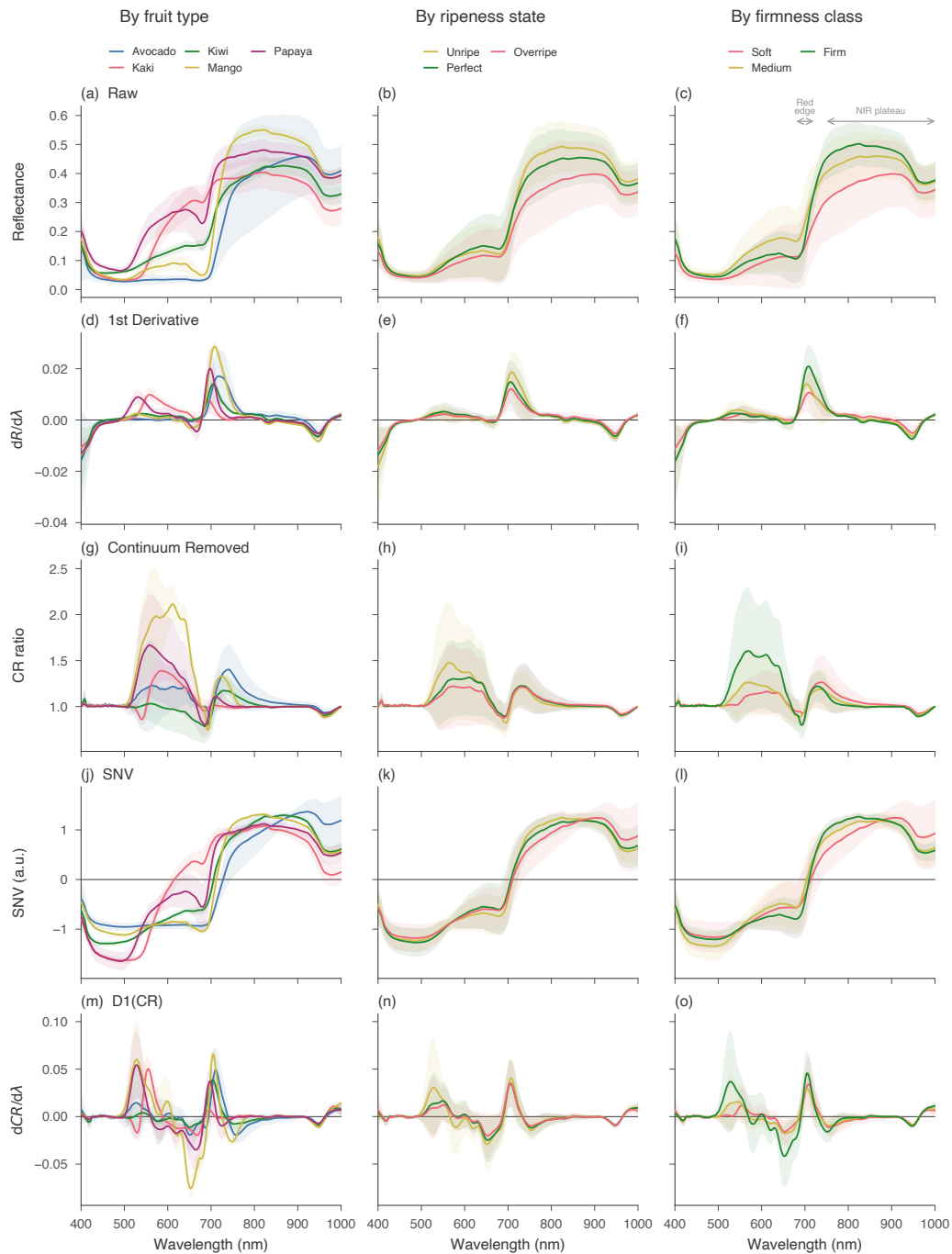


Figure 2: Five spectral preprocessing types (rows) across the 224 bands (400–1000 nm), grouped by fruit type, ripeness state, and firmness class (columns). Shaded bands: mean  $\pm$  std. Rows: (a–c) raw, (d–f) 1st derivative, (g–i) continuum removed, (j–l) SNV, (m–o) D1(CR). Inter-species variation dominates; class separation is clearest in derivative features near the red edge ( $\sim 700$  nm).

227 keeping the benchmark test set fixed. This changes both the sample compo-  
228 sition and size (an 8.7% increase in training samples), so its advantage may  
229 partly reflect additional data rather than purely improved balance.

230 *Training set augmentation.* SMOTE, random oversampling, and random  
231 undersampling modify the class distribution within the original training par-  
232 tition without changing which real samples are included. Sampling is applied  
233 independently per task (ripeness or firmness); the test partition is kept fixed.

234 Firmness is discretized into three bins: soft (0–1000 gf), medium (1001–  
235 2500 gf), and firm (2501+ gf), derived from cross-fruit firmness distribution  
236 analysis. The original DeepHS dataset defines fruit-specific thresholds (e.g.,  
237  $<900/900\text{--}1200/>1200$  g/cm<sup>2</sup> for avocado; Varga et al. 2021), but universal  
238 bins enable direct cross-fruit comparison across the full firmness spectrum.

239 The original unbalanced split is preserved for comparability with bench-  
240 marks from Varga et al. (Varga et al., 2021) and Ben Jmaa et al. (Ben Jmaa  
241 et al., 2025).

#### 242 3.2.4. Additional Preprocessing Strategies

243 When enabled, PCA reduces dimensionality while retaining 95% of vari-  
244 ance, compressing the 1,120-dimensional feature space to approximately 18–  
245 20 principal components. PCA is integrated within the scikit-learn Pipeline  
246 class: the transformation is fit only on training data and applied to test data,  
247 and in Phase 4 cross-validation PCA is fit separately within each fold. The  
248 five spectral representations are concatenated before PCA, so it operates on  
249 the full combined feature space.

250 SMOTE (Synthetic Minority Oversampling Technique) generates syn-

251 thetic samples for minority classes using the imbalanced-learn<sup>2</sup> library ( $k=5$   
252 nearest neighbors), applied independently per task after train/test splitting  
253 but before model training. For the firmness task, which includes 14 samples  
254 (2.7%) with missing values encoded as a separate class, SMOTE treats this as  
255 a fourth firmness category. Random oversampling duplicates minority class  
256 samples; random undersampling removes majority class samples.

257 Phase 2 evaluates these strategies across 200 configurations (5 data bal-  
258 ancing strategies  $\times$  2 PCA options  $\times$  20 models).

### 259 3.2.5. Wavelength Subset Experiments (*VIS-3* and *RGB*)

260 Two three-band configurations are evaluated to assess the necessity of  
261 full-spectrum sensing.

262 *VIS-3 (XAI-derived)*. Three visible-range bands are selected from hyper-  
263 spectral data at indices {18, 52, 89}, corresponding to wavelengths {448, 540, 640} nm.  
264 These bands were identified through consensus feature ranking from Ex-  
265 traTrees explainability analysis (Section 4), selecting the top joint-ranked  
266 visible-range bands (<700 nm) across both ripeness and firmness tasks. We  
267 term this the “VIS-3 subset” to emphasize that these bands are selected from  
268 hyperspectral cubes, not acquired by an independent RGB sensor.

269 *RGB (approximate)*. Three bands at indices {19, 56, 93}, corresponding  
270 to wavelengths {450, 550, 651} nm, approximate the center wavelengths of  
271 standard Bayer-pattern RGB sensors. This configuration tests whether XAI-  
272 guided band selection outperforms generic RGB-center wavelengths, thereby  
273 isolating the contribution of precise band placement from the effect of spectral

---

<sup>2</sup>imbalanced-learn: <https://imbalanced-learn.org>

274 dimensionality reduction.

275 Both configurations apply the same five spectral transformations, result-  
276 ing in a 15-dimensional feature vector (3 bands  $\times$  5 methods) compared to  
277 the 1,120-dimensional full-spectrum vector. Derivatives and continuum re-  
278 moval of three non-contiguous bands lack the physical meaning they carry  
279 over continuous spectra; reduced-band results therefore represent an upper  
280 bound on what a true 3-band sensor could achieve, not a realistic sensor  
281 comparison. Both pipelines follow identical experimental procedures to the  
282 full-spectrum pipeline.

### 283 *3.3. Experimental Design*

284 The six-phase methodology trains each algorithm separately for ripeness  
285 and firmness, with performance reported jointly through aggregated metrics.

#### 286 *3.3.1. Pipeline Phases*

287 Each model is trained separately for ripeness and firmness prediction,  
288 with performance evaluated using overall accuracy (OA, Eq. 11). We term  
289 this “paired single-task evaluation” to distinguish it from multi-task learning  
290 architectures that share representations across tasks. Each task is modeled  
291 independently; only the evaluation metric aggregates performance. We use  
292 separate models rather than joint multi-output classification for four rea-  
293 sons: (1) many evaluated algorithms do not natively support multi-output  
294 classification; (2) separate models allow independent hyperparameter tuning  
295 (firmness accuracy exceeds ripeness accuracy by approximately 25 pp, Ta-  
296 ble 5); (3) agricultural operations may prioritize one attribute depending on  
297 market demands; and (4) existing literature predominantly uses single-task

298 evaluation. This approach does not exploit shared representations that joint  
299 models might capture; future work should compare separate-task versus joint  
300 multi-task architectures.

301 Phase 1 establishes baseline performance for 20 algorithms on the orig-  
302 inal imbalanced dataset using default hyperparameters, evaluated using ac-  
303 curacy, F1-macro, and F1-weighted metrics for each task. Phase 2 evaluates  
304 preprocessing strategies using a full-factorial design across data balancing  
305 strategy (five levels: original imbalanced, stratified resplit, SMOTE, ran-  
306 dom oversampling, random undersampling), PCA dimensionality reduction  
307 (enabled at 95% variance or disabled), and model selection (20 algorithms),  
308 generating 200 configurations. The best Phase 2 preprocessing setting per  
309 modality is selected for Phase 3 based on overall accuracy. This selection  
310 uses test-set performance, introducing indirect information leakage: the cho-  
311 sen preprocessing was optimized, in part, on the same held-out set used for  
312 final evaluation. Phase 4 cross-validation provides a partial mitigation, but  
313 the preprocessing choice itself was not cross-validated.

314 Phase 3 applies Bayesian optimization via Optuna<sup>3</sup> (Akiba et al., 2019)  
315 with the best Phase 2 preprocessing configuration. For each model, 100  
316 trials using Tree-structured Parzen Estimator (TPE) sampling (Bergstra  
317 et al., 2011) maximize  $\overline{F1}_{\text{macro}}$  (Eq. 12). Search spaces are model-specific  
318 (Table B.1): tree-based models optimize depth, minimum samples per leaf,  
319 and ensemble size; linear models optimize regularization; neural networks op-  
320 timize learning rate, hidden layer sizes, and activations. Optimization uses

---

<sup>3</sup>Optuna: <https://optuna.org>

321 stratified 70/30 train/validation splits, with final models retrained on the full  
322 training set. Total: 2,000 trials per modality (20 models  $\times$  100 trials), 4,000  
323 trials across both modalities.

324 Phase 4 validates model performance through 10-fold stratified cross-  
325 validation using optimal hyperparameters from Phase 3 and best preprocess-  
326 ing from Phase 2. Results are reported as mean  $\pm$  standard deviation across  
327 folds, with 95% confidence intervals and coefficient of variation. Phases 1–3  
328 provide exploratory results on a single fixed split; Phase 4 cross-validation is  
329 the authoritative basis for algorithm ranking and all primary conclusions.

330 Phase 5 combines top models from Phase 4 using soft voting, hard voting,  
331 stacking, and blending. Phase 6 performs dedicated explainability analysis  
332 for the best-performing single model using SHAP<sup>4</sup> (Lundberg and Lee, 2017),  
333 LIME<sup>5</sup> (Ribeiro et al., 2016), and permutation importance (Breiman, 2001),  
334 producing wavelength-level interpretations for both tasks. The full pipeline  
335 architecture is illustrated in Fig. 3.

### 336 3.3.2. Benchmark Reference

337 We reference Fruit-HSNet’s reported 70.73% overall accuracy (Ben Jmaa  
338 et al., 2025) as a contextual data point, not a controlled comparison. Fruit-  
339 HSNet processes full 64 $\times$ 64 pixel hyperspectral cubes through a dual-branch  
340 network that extracts Fourier-domain spectral features alongside central-  
341 pixel spatial signatures, trained with GPU acceleration, focal loss, and spatial  
342 data augmentation including random cropping, flipping, rotation, and noise

---

<sup>4</sup>SHAP: <https://github.com/shap/shap>

<sup>5</sup>LIME: <https://github.com/marcotcr/lime>

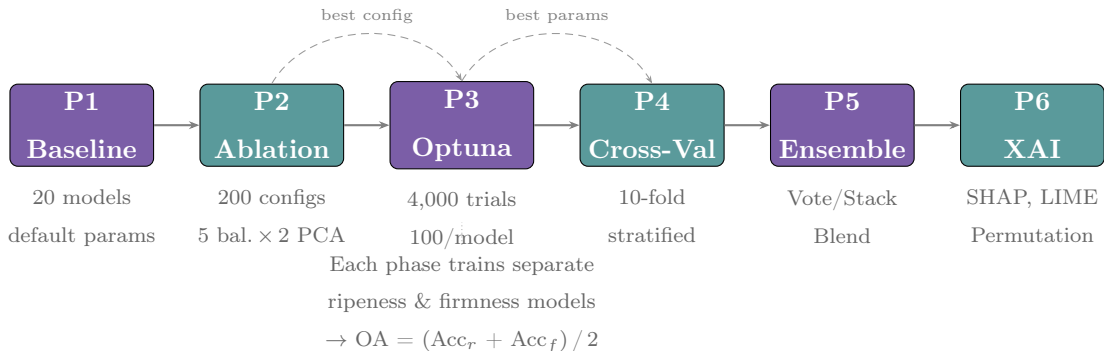


Figure 3: Six-phase benchmark pipeline. Each phase trains independent ripeness and firmness models, combined via overall accuracy (OA). Dashed arrows indicate configuration propagation between phases.

343 injection. Our approach discards spatial information entirely, representing  
 344 each sample by its per-pixel-averaged mean spectrum (a 224-dimensional vec-  
 345 tor retaining no spatial structure) processed on a consumer CPU in under  
 346 one second. The following protocol differences preclude direct comparison:  
 347 (1) Fruit-HSNet used a different train/test partition, (2) label preprocess-  
 348 ing and class definitions may differ, and (3) Fruit-HSNet is a deep learn-  
 349 ing architecture while this study evaluates the scikit-learn ecosystem. Our  
 350 Phase 3 result (ExtraTrees, 75.0%) carries a Wilson 95% confidence inter-  
 351 val of [67.6%, 81.8%], which contains Fruit-HSNet’s reported value, and no  
 352 statistically meaningful difference can be claimed.

353 Phase 4 cross-validation provides statistical estimates through 10-fold  
 354 stratified sampling (mean, standard deviation, coefficient of variation, 95%  
 355 confidence intervals). We apply the Friedman test for omnibus comparison  
 356 of all 20 classifiers on per-fold F1 scores, followed by Nemenyi post-hoc tests  
 357 with family-wise error control (Table 2). Cohen’s  $d$  effect sizes quantify task

358 asymmetry (Table 3). Phase 4 cross-validated rankings are the authoritative  
359 ranking for algorithm comparison; Phase 3 single-split results are retained  
360 only for protocol-matched reference against Fruit-HSNet. Reproducibility  
361 details are provided in Table B.3.

### 362 3.4. Machine Learning Models

363 We evaluate 20 classical and gradient-boosted machine learning algo-  
364 rithms from seven algorithmic families, including PLS-DA (Partial Least  
365 Squares Discriminant Analysis), the standard chemometric baseline for spec-  
366 tral classification. PLS-DA is implemented as a custom scikit-learn wrapper  
367 around `PLSRegression` with dummy-coded class labels and `argmax` predic-  
368 tion; `n_components` is clamped at fit time to  $\min(n_{\text{samples}}, n_{\text{features}}, n_{\text{classes}})$ .  
369 Note that the sklearn `MLPClassifier` used here is a shallow multi-layer per-  
370 ceptron; it shares an API with other scikit-learn classifiers but is not a deep  
371 learning model. Models are implemented using scikit-learn<sup>6</sup> (Pedregosa et al.,  
372 2011), XGBoost<sup>7</sup> (Chen and Guestrin, 2016), and LightGBM<sup>8</sup> (Ke et al.,  
373 2017). Each model is integrated into a pipeline with feature scaling (Stan-  
374 dardScaler) and optional PCA. Table B.1 provides a structured overview  
375 of all models, including abbreviations, family membership, hyperparameter  
376 counts, Optuna search spaces, and native feature importance support.

377 All models use fixed random state (`seed=42`), early stopping where appli-  
378 cable, and class balancing for imbalanced datasets. A complete overview of  
379 all models with Optuna search spaces is provided in Table B.1 (Appendix);

---

<sup>6</sup>scikit-learn: <https://scikit-learn.org>

<sup>7</sup>XGBoost: <https://xgboost.ai>

<sup>8</sup>LightGBM: <https://lightgbm.readthedocs.io>

380 best hyperparameters from Bayesian optimization are listed in Table B.2.

### 381 3.5. Evaluation Metrics

382 We report per-task accuracies for ripeness ( $r$ ) and firmness ( $f$ ), overall  
 383 accuracy (OA) defined as the arithmetic mean of task accuracies, and F1-  
 384 scores (both macro-averaged and weighted). Per-class metrics for class  $c$ :

$$\text{Precision}_c = \frac{TP_c}{TP_c + FP_c}, \quad (6)$$

$$\text{Recall}_c = \frac{TP_c}{TP_c + FN_c}, \quad (7)$$

$$\text{F1}_c = \frac{2 \text{Precision}_c \text{Recall}_c}{\text{Precision}_c + \text{Recall}_c}. \quad (8)$$

387 Task-level aggregates for task  $t \in \{r, f\}$  with class set  $C_t$ :

$$\text{F1}_{\text{macro}}^{(t)} = \frac{1}{|C_t|} \sum_{c \in C_t} \text{F1}_c, \quad (9)$$

$$\text{F1}_{\text{wt}}^{(t)} = \sum_{c \in C_t} \frac{n_c}{N_t} \text{F1}_c. \quad (10)$$

389 Overall metrics average across both tasks:

$$\text{OA} = \frac{\text{Acc}^{(r)} + \text{Acc}^{(f)}}{2}, \quad (11)$$

$$\overline{\text{F1}}_{\text{macro}} = \frac{\text{F1}_{\text{macro}}^{(r)} + \text{F1}_{\text{macro}}^{(f)}}{2}. \quad (12)$$

391 Equal weighting treats both quality attributes as equally relevant for  
 392 holistic assessment; deployment scenarios prioritizing one attribute can reweight  
 393 accordingly. OA is the primary benchmark comparison metric;  $\overline{\text{F1}}_{\text{macro}}$  is  
 394 the optimization objective for Phase 3 and Phase 4 cross-validation. Cross-  
 395 validation reports mean  $\pm$  standard deviation, coefficient of variation, and  
 396 95% confidence intervals. Efficiency metrics include training time (s), infer-  
 397 ence time (s), and model size (MB).

398 *3.6. Computational Environment*

399 All experiments were conducted on consumer-grade hardware, confirm-  
400 ing feasibility on consumer hardware without specialized infrastructure (Ta-  
401 ble B.3). Hardware configuration consists of MacBook Air M4 (2024) with  
402 Apple Silicon M4 chip, 24 GB unified memory, SSD storage, and macOS  
403 (arm64). Software stack includes Python 3.13, scikit-learn 1.7.1, XGBoost  
404 3.0.4, LightGBM 4.6.0, Optuna 4.5.0, and imbalanced-learn 0.14.0.

405 CO<sub>2</sub> emissions were negligible across all models (<0.1 g per model on  
406 CPU-only consumer hardware; estimated via CodeCarbon<sup>9</sup> methodology (Schmidt  
407 et al., 2021)).

408 **4. Results**

409 Evaluation spans 20 algorithms across three modalities (full-spectrum:  
410 1,120 features; VIS-3 subset: 15 features; RGB approximate: 15 features)  
411 and six experimental phases, including 10-fold cross-validation and 4,000  
412 hyperparameter-optimization trials. A Friedman test on per-fold overall F1  
413 scores confirmed significant differences among classifiers ( $\chi^2(19) = 144.8$ ,  
414  $p < 10^{-20}$ ); the Nemenyi post-hoc critical difference diagram (Fig. A.1 in  
415 Appendix, Table 2) identifies 46 significantly different pairs.

416 *4.1. Pipeline Performance*

417 Table 1 presents integrated results across all pipeline phases for the 20  
418 evaluated models, sorted by Phase 4 cross-validated overall accuracy. XG-  
419 Boost achieved the highest cross-validated OA ( $82.3 \pm 5.4\%$ ), with HistGra-

---

<sup>9</sup>CodeCarbon: <https://codecarbon.io>

420 dientBoosting, LGBM, and ExtraTrees within 0.1 percentage points (82.2%  
421 each). In total, 11 of 20 models exceeded 75% OA in 10-fold CV. ExtraTrees,  
422 which ranked first on the single benchmark split (Phase 3, 75.0%), maintained  
423 strong cross-validated performance. Phase 4 cross-validated rankings are the  
424 authoritative comparison. Subsequent per-class and XAI analyses use Ex-  
425 traTrees because it provides native feature importance scores and achieved  
426 the highest single-split accuracy on the benchmark partition used for Fruit-  
427 HSNNet comparison. Training time was 0.2 seconds on consumer hardware.

428 Preprocessing strategy substantially influenced performance for top mod-  
429 els. Stratified resplit balancing improved ExtraTrees from 71.7% (Phase 1  
430 baseline) to 75.0% (Phase 2), a 3.3 percentage point gain from data bal-  
431 ancing alone. Phase 3 Bayesian optimization (100 trials per model) did not  
432 further improve ExtraTrees. For competitive models, preprocessing rather  
433 than hyperparameter tuning drove the performance gain (Fig. 4a). All 20  
434 models trained in under 19 seconds on the consumer-grade hardware (Mac-  
435 Book Air M4). Inference latency (Phase 4) is well below the 100 ms sorting-  
436 line threshold for all top models: ExtraTrees 2.1 ms/sample, RandomForest  
437 2.9 ms/sample, LGBM 6.8 ms/sample, XGBoost 34.1 ms/sample (138 test  
438 samples). The most efficient model (ExtraTrees) is 90× faster than Gradi-  
439 entBoosting while achieving 11.2 percentage points higher accuracy (Fig. 4b).

440

441 The Pareto frontier in model size vs. accuracy space (Fig. 4b, star mark-  
442 ers) identifies six non-dominated models, namely ExtraTrees (75.0%, 2.8 MB),  
443 RandomForest (73.6%, 1.4 MB), XGBoost (65.6%, 0.7 MB), GradientBoost-  
444 ing (63.8%, 0.4 MB), AdaBoost (62.0%, 0.1 MB), and DecisionTree (58.7%,

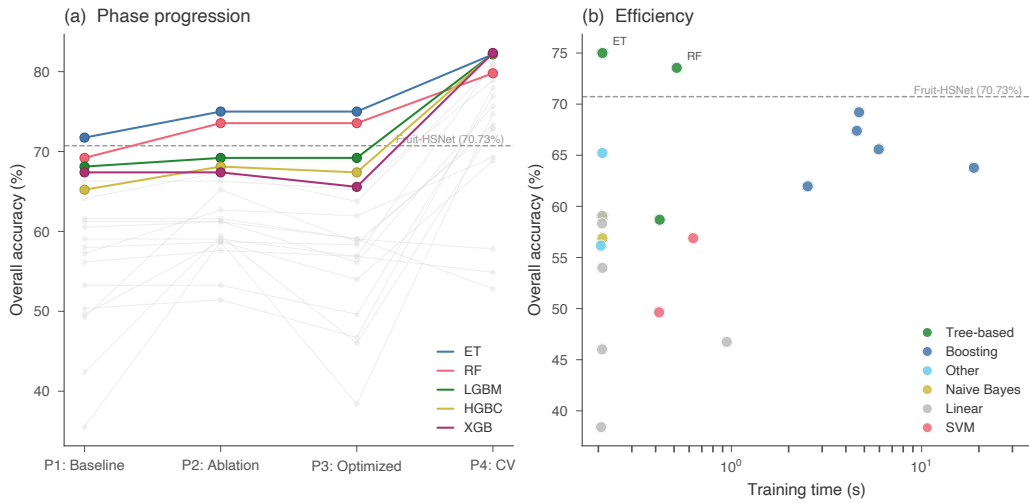


Figure 4: Pipeline performance. (a) Phase progression for all 20 models: connected dots trace each model’s overall accuracy from baseline (P1) through ablation (P2, best config), optimization (P3), and 10-fold cross-validation (P4). Top-5 models highlighted; remaining in gray. The largest gains occur between P1 and P2, confirming that preprocessing contributes more than hyperparameter tuning for competitive models. (b) Training-time efficiency (log scale) vs overall accuracy (Phase 3), colored by algorithmic family. Marker size reflects model storage size. Star markers denote Pareto-optimal models in model size vs accuracy space (6 models: ExtraTrees, RandomForest, XGBoost, GradientBoosting, AdaBoost, DecisionTree), representing the best size-accuracy tradeoff frontier. Tree-based methods achieve the highest accuracy with the shortest training times.

445 0.07 MB). ExtraTrees achieves the highest accuracy at a storage footprint  
446 suitable for embedded sorting hardware, with inference latency of 2.1 ms per  
447 sample. The dominance of preprocessing over optimization across competi-  
448 tive models raises the question of which preprocessing decisions matter most  
449 and why.

Table 1: Model performance and efficiency across pipeline phases (full-spectrum, 1,120 features). Models sorted by Phase 4 cross-validated overall accuracy.<sup>11</sup> P4 shows 10-fold stratified cross-validation accuracy (mean  $\pm$  std). P3 Wilson 95% CI (binomial,  $n=138$ ) shown in brackets. **Bold** = best in column; underline = second-best. Arrows indicate preferred direction.

Rk	Model	Family	P1 $\uparrow$	P2 $\uparrow$	P2 Config	P3 $\uparrow$	P3 CI	P4 CV (%) $\uparrow$	Time $\downarrow$	Size $\downarrow$
			OA	OA		OA	[95%]		(s)	(MB)
1	XGBoost	Boosting	67.4	67.4	SMOTE / No PCA	65.6	[57.0, 72.7]	<b>82.3</b> $\pm$ 5.4	5.9	0.72
2	HistGradientBoosting	Boosting	65.2	68.1	SMOTE / No PCA	67.4	[59.2, 74.6]	<u>82.2</u> $\pm$ 5.4	4.6	5.05
3	LGBM	Boosting	68.1	69.2	Strat. Resplit / No PCA	69.2	[61.4, 76.6]	<u>82.2</u> $\pm$ 4.9	4.7	1.64
4	ExtraTrees	Tree-based	<b>71.7</b>	<b>75.0</b>	Strat. Resplit / No PCA	<b>75.0</b>	[67.6, 81.8]	<u>82.2</u> $\pm$ 4.7	<b>0.2</b>	2.77
5	SVC_RBF	SVM	56.2	57.6	SMOTE / No PCA	56.9	[48.2, 64.5]	81.1 $\pm$ 4.0	0.6	5.16
6	GradientBoosting	Boosting	64.1	67.4	Oversample / No PCA	63.8	[55.5, 71.3]	80.6 $\pm$ 5.8	18.8	0.41
7	RandomForest	Tree-based	<u>69.2</u>	<u>73.6</u>	Strat. Resplit / No PCA	<u>73.6</u>	[66.0, 80.5]	79.8 $\pm$ 4.8	0.5	1.44
8	KNeighbors	Other	66.3	66.3	Strat. Resplit / PCA	65.2	[57.0, 72.7]	78.9 $\pm$ 4.4	<b>0.2</b>	7.14
9	LogisticRegression	Linear	50.4	51.4	Oversample / No PCA	46.7	[38.3, 54.7]	78.0 $\pm$ 5.2	0.9	0.11
10	SVC_Linear	SVM	53.3	53.3	Original / No PCA	49.6	[41.1, 57.5]	76.9 $\pm$ 5.2	<u>0.4</u>	3.64
11	MLPClassifier	Other	61.2	61.2	Original / No PCA	56.2	[48.2, 64.5]	75.6 $\pm$ 6.0	<b>0.2</b>	5.22
12	LDA	Linear	35.5	58.7	Strat. Resplit / PCA	38.4	[30.7, 46.7]	74.7 $\pm$ 4.6	<b>0.2</b>	0.23
13	SGDClassifier	Linear	58.0	58.7	Strat. Resplit / PCA	58.3	[49.6, 65.9]	73.4 $\pm$ 4.2	<b>0.2</b>	0.11
14	DecisionTree	Tree-based	49.3	65.2	SMOTE / No PCA	58.7	[50.4, 66.6]	72.8 $\pm$ 5.2	<u>0.4</u>	0.07
15	Ridge	Linear	42.4	59.4	Strat. Resplit / PCA	46.0	[38.3, 54.7]	72.8 $\pm$ 3.9	<b>0.2</b>	0.11
16	AdaBoost	Boosting	57.2	62.7	Original / PCA	62.0	[54.0, 70.0]	69.3 $\pm$ 5.9	2.5	0.12
17	PassiveAggressive	Linear	49.6	58.7	SMOTE / PCA	54.0	[45.3, 61.7]	68.8 $\pm$ 5.4	<b>0.2</b>	0.11
18	PLSDA	Chemometric	61.6	61.6	Original / No PCA	59.1	[51.1, 67.3]	57.8 $\pm$ 4.5	<b>0.2</b>	0.42
19	BernoulliNB	Naive Bayes	59.1	59.1	Undersample / No PCA	56.9	[48.2, 64.5]	54.9 $\pm$ 5.5	<b>0.2</b>	0.17
20	GaussianNB	Naive Bayes	60.5	61.2	SMOTE / PCA	59.1	[51.1, 67.3]	52.8 $\pm$ 4.9	<b>0.2</b>	0.17

Table 2: Model rankings from 10-fold cross-validation (Friedman/Nemenyi,  $\alpha = 0.05$ ). Rank = mean rank across folds (lower is better). Models connected by a bar in Fig. A.1 are not significantly different.

Rank	Model	Avg. Rank↓	Mean F1↑	Family
1	ET	<b>3.7</b>	<b>0.792 ± 0.077</b>	Tree-based
2	SVM-R	<u>4.7</u>	<u>0.787 ± 0.077</u>	SVM
3	KNN	5.4	0.771 ± 0.061	Other
4	GBC	5.8	0.760 ± 0.064	Boosting
5	RF	6.2	0.758 ± 0.063	Tree-based
6	LGBM	6.2	0.755 ± 0.049	Boosting
7	HGBC	6.9	0.752 ± 0.048	Boosting
8	SVM-L	7.2	0.747 ± 0.067	SVM
9	LR	7.5	0.750 ± 0.077	Linear
10	XGB	7.6	0.743 ± 0.063	Boosting
11	MLP	7.7	0.725 ± 0.049	Other
12	LDA	12.2	0.671 ± 0.060	Linear
13	DT	13.0	0.659 ± 0.052	Tree-based
14	SGD	14.0	0.639 ± 0.055	Linear
15	Ridge	14.6	0.626 ± 0.041	Linear
16	PA	15.5	0.610 ± 0.057	Linear
17	Ada	15.7	0.609 ± 0.044	Boosting
18	BNB	18.1	0.526 ± 0.065	Naive Bayes
19	PLS-DA	18.7	0.491 ± 0.065	Chemometric
20	GNB	19.3	0.486 ± 0.049	Naive Bayes

Table 3: Summary of statistical tests and descriptive comparisons. Effect sizes follow Cohen’s conventions:  $|d| < 0.2$  negligible, 0.2–0.5 small, 0.5–0.8 medium,  $> 0.8$  large. Friedman and Nemenyi tests are appropriate because models share the same CV folds. Firmness vs. ripeness comparisons involve the same 20 models trained on shared data; formal significance tests are not used because the observations are not independent. Phase 3 modality comparisons (HSI vs. VIS-3, HSI vs. RGB, VIS-3 vs. RGB) use paired Wilcoxon signed-rank tests on per-model overall accuracy.

Comparison	Summary	$d$	Effect
Friedman (20 classifiers)	$\chi^2(19) = 144.8, p < 10^{-20}$	—	—
Nemenyi pairs ( $\alpha=0.05$ )	46 of 190 pairs differ	—	—
Firmness > Ripeness (IQR: 21.4–32.8 pp; exception: LDA)	19/20 models; mean gap 25.0 pp	1.84	Large
HSI vs. RGB	Cohen’s $d$ from Phase 4 fold F1	1.17	Large
HSI vs. VIS-3	Phase 3 OA, 20 models; $p = 0.46$	0.32	Small
HSI vs. RGB	Phase 3 OA, 20 models; $p = 0.39$	0.14	Negl.
VIS-3 vs. RGB	Phase 3 OA, 20 models; $p = 0.30$	0.27	Small

450 *4.2. Preprocessing Ablation*

451 Table 4 summarizes the Phase 2 ablation study across 200 configurations  
452 (5 balancing strategies  $\times$  2 PCA options  $\times$  20 models). Stratified resplit  
453 achieved the highest mean overall accuracy (59.08% without PCA), while  
454 random undersampling performed worst (46.56%).

455 PCA degraded performance in four of five balancing conditions. Across  
456 all 200 configurations, models without PCA outperformed those with PCA  
457 by 2.46 percentage points on average (55.62% vs 53.16%, Fig. 5). Only 7  
458 of 20 models utilized PCA in their best Phase 2 configuration, and these  
459 were predominantly linear models (Ridge, LDA, SGDClassifier) rather than  
460 tree-based methods. The full 200-configuration landscape is visualized in  
461 Fig. 5, where the stratified resplit / no PCA column consistently achieves  
462 the warmest values.

463 The consistent PCA degradation is surprising given its routine use in  
464 hyperspectral analysis and warrants mechanistic explanation.

465 Across the  $20 \times 10$  Phase 2 accuracy matrix, the median accuracy range  
466 attributable to preprocessing (varying balancing strategy and PCA within  
467 each model) is 25.2 pp, compared with 32.4 pp for algorithm selection (vary-  
468 ing the model within each preprocessing configuration). A bootstrap analysis  
469 (10,000 resamples) yields a 95% CI of  $[-5.8, 21.9]$  pp on the difference be-  
470 tween these median ranges, an interval that contains zero. Preprocessing  
471 and algorithm choice are thus of comparable magnitude on this benchmark.  
472 Ridge exemplifies the upper bound of preprocessing impact: its overall accu-  
473 racy rose from 42.4% (Phase 1 baseline) to 59.4% (Phase 2 best), a 17.0 pp  
474 gain from data engineering alone.

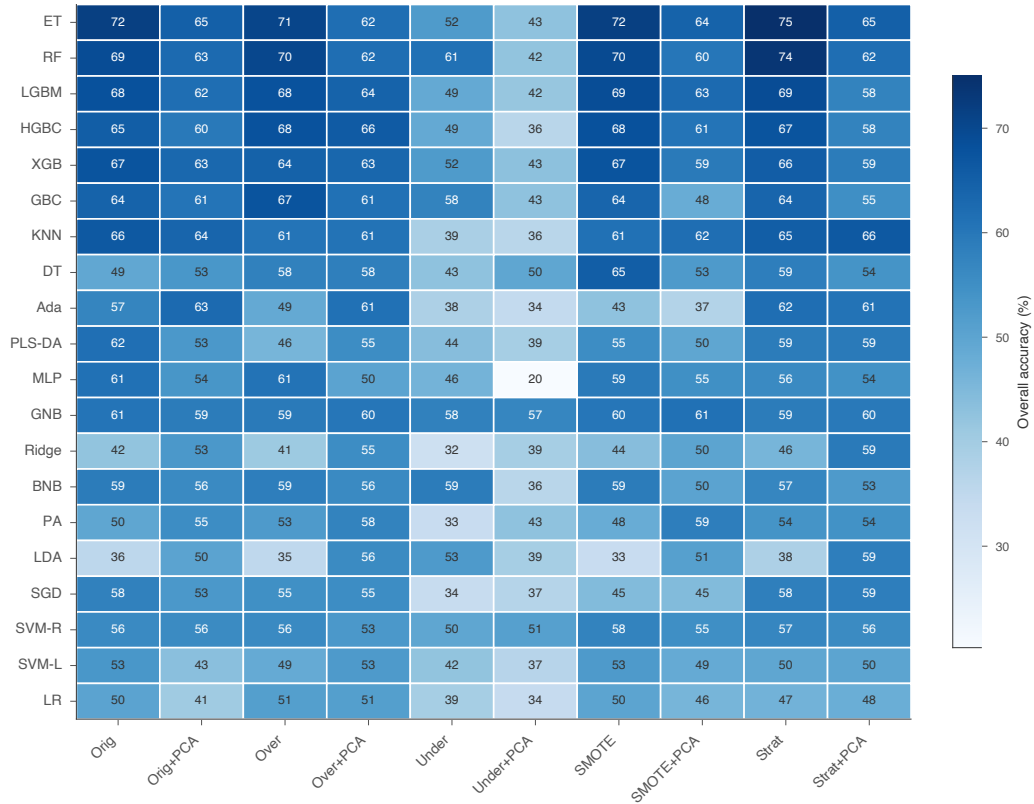


Figure 5: Phase 2 preprocessing ablation across 200 configurations. Overall accuracy (%) for 20 models (rows, sorted by best OA) across 10 preprocessing configurations (5 balancing strategies  $\times$  2 PCA settings). Vertical white lines separate balancing strategies; within each pair, the left column is without PCA and the right is with PCA. The stratified resplit / no PCA column consistently achieves the highest values.

Table 4: Phase 2 ablation: mean  $\pm$  SD overall accuracy (%) $\uparrow$  across 20 models by balancing strategy and PCA setting.  $\Delta$  = no-PCA minus PCA difference. Best strategy in bold. Large SDs (4–11 pp) reflect substantial model-to-model variance within each configuration.  $\dagger$ Stratified resplit replaces the training set ( $n=414$  vs  $n=381$ ); strategies below the rule augment the original training set.

Balancing strategy	No PCA	PCA	$\Delta$ (pp)
<b>Stratified resplit<math>\dagger</math></b>	<b>59.08<math>\pm</math>9.26</b>	<b>57.41<math>\pm</math>4.57</b>	+1.67
SMOTE	57.14 $\pm$ 10.61	53.84 $\pm$ 7.18	+3.30
Original unbalanced	58.32 $\pm$ 9.38	56.29 $\pm$ 6.59	+2.03
Random oversampling	57.01 $\pm$ 9.85	58.12 $\pm$ 4.47	-1.11
Random undersampling	46.56 $\pm$ 8.93	40.13 $\pm$ 7.47	+6.43
<i>Overall mean</i>	55.62	53.16	+2.46

#### 4.3. Task Asymmetry and Modality Comparison

Table 5 and Fig. 6 reveal two consistent patterns across all 20 models: (1) firmness prediction outperforms ripeness by a wide margin (Fig. 6a), and (2) full-spectrum data outperforms VIS-3 and RGB subsets for most competitive models, but the advantage is small (Fig. 6b).

Firmness accuracy exceeded ripeness accuracy in 19 of 20 models (mean gap 25.0 percentage points; IQR: 21.4–32.8 pp; Cohen’s  $d = 1.84$ , large effect). This asymmetry was most extreme for SVC\_RBF (31.9% ripeness vs 81.9% firmness, a 50.0 pp gap) and smallest for SVC\_Linear (44.9% vs 54.3%, a 9.4 pp gap). LDA was the only model where ripeness accuracy exceeded firmness, likely reflecting model-specific failure modes. Because the 20 models share a common training set and test set, formal significance tests

487 for this comparison are not appropriate (Table 3).

488 Phase 4 fold-level F1 comparisons between HSI and RGB revealed a large  
489 full-spectrum advantage (Cohen’s  $d = 1.17$ ; Table 3), concentrated among  
490 the top-performing models. In Phase 3, VIS-3 and RGB yield compara-  
491 ble performance (mean OA 61.7% vs 60.1%; Cohen’s  $d = 0.27$ ,  $p = 0.30$ ),  
492 with neither modality significantly outperforming the other (Table 3). Both  
493 reduced-band configurations outperformed full-spectrum on average across  
494 all 20 models (61.7% and 60.1% vs 59.1%), but this average is driven by  
495 weak models (Ridge, LDA, and SVC\_Linear) that perform substantially  
496 better in the low-dimensional 15-feature space, likely because dimensionality  
497 reduction mitigates the multicollinearity that destabilizes their coefficient es-  
498 timates. At the top of the leaderboard, the full-spectrum advantage persists:  
499 ExtraTrees achieves 75.0% (full) vs 72.1% (RGB) vs 71.0% (VIS-3). The  
500 persistence of the firmness–ripeness gap across all 20 models and all three  
501 modalities suggests a property of the labels rather than of the algorithms.

Table 5: Task asymmetry and modality comparison of Phase 3 models. RA = ripeness accuracy, FA = firmness accuracy, OA = overall accuracy. Full-spectrum has 1,120 features; VIS-3 and RGB each use 15 features.  $\Delta\text{OA} = \text{Full OA} - \max(\text{VIS-3}, \text{RGB})$ .

Rk	Model	Full-spectrum			Reduced-band OA $\uparrow$		
		RA $\uparrow$	FA $\uparrow$	OA $\uparrow$	VIS-3	RGB	$\Delta\text{OA}$
1	ExtraTrees	<b>63.0</b>	<b>87.0</b>	<b>75.0</b>	<b>71.0</b>	<b>72.1</b>	+2.9
2	RandomForest	<u>60.9</u>	<u>86.2</u>	<u>73.6</u>	67.8	<u>71.4</u>	+2.2
3	LGBM	58.0	80.4	69.2	68.1	68.8	+0.4
4	HistGradBoosting	54.3	80.4	67.4	67.4	69.9	-2.5
5	XGBoost	52.2	79.0	65.6	<u>68.5</u>	68.5	-2.9
6	KNeighbors	48.6	81.9	65.2	65.6	66.7	-1.4
7	GradientBoosting	47.1	80.4	63.8	67.4	65.9	-3.6
8	AdaBoost	48.6	75.4	62.0	60.9	44.9	+1.1
9	GaussianNB	44.2	73.9	59.1	54.0	52.5	+5.1
10	PLSDA	39.9	78.3	59.1	51.8	48.6	+7.2
11	DecisionTree	50.7	66.7	58.7	58.3	60.9	-2.2
12	SGDClassifier	42.0	74.6	58.3	60.9	57.2	-2.5
13	BernoulliNB	41.3	72.5	56.9	60.5	53.3	-3.6
14	SVC_RBF	31.9	81.9	56.9	61.2	58.7	-4.3
15	MLPClassifier	39.1	73.2	56.2	51.8	58.0	-1.8
16	PassiveAggr.	42.0	65.9	54.0	58.7	60.1	-6.2
17	SVC_Linear	44.9	54.3	49.6	64.1	57.2	-14.5
18	LogisticRegr.	37.7	55.8	46.7	47.5	56.5	-9.8
19	Ridge	37.0	55.1	46.0	66.7	56.5	-20.7
20	LDA	47.8	29.0	38.4	61.2	54.7	-22.8
<i>Mean</i>		46.6	71.6	59.1	61.7	60.1	-4.0

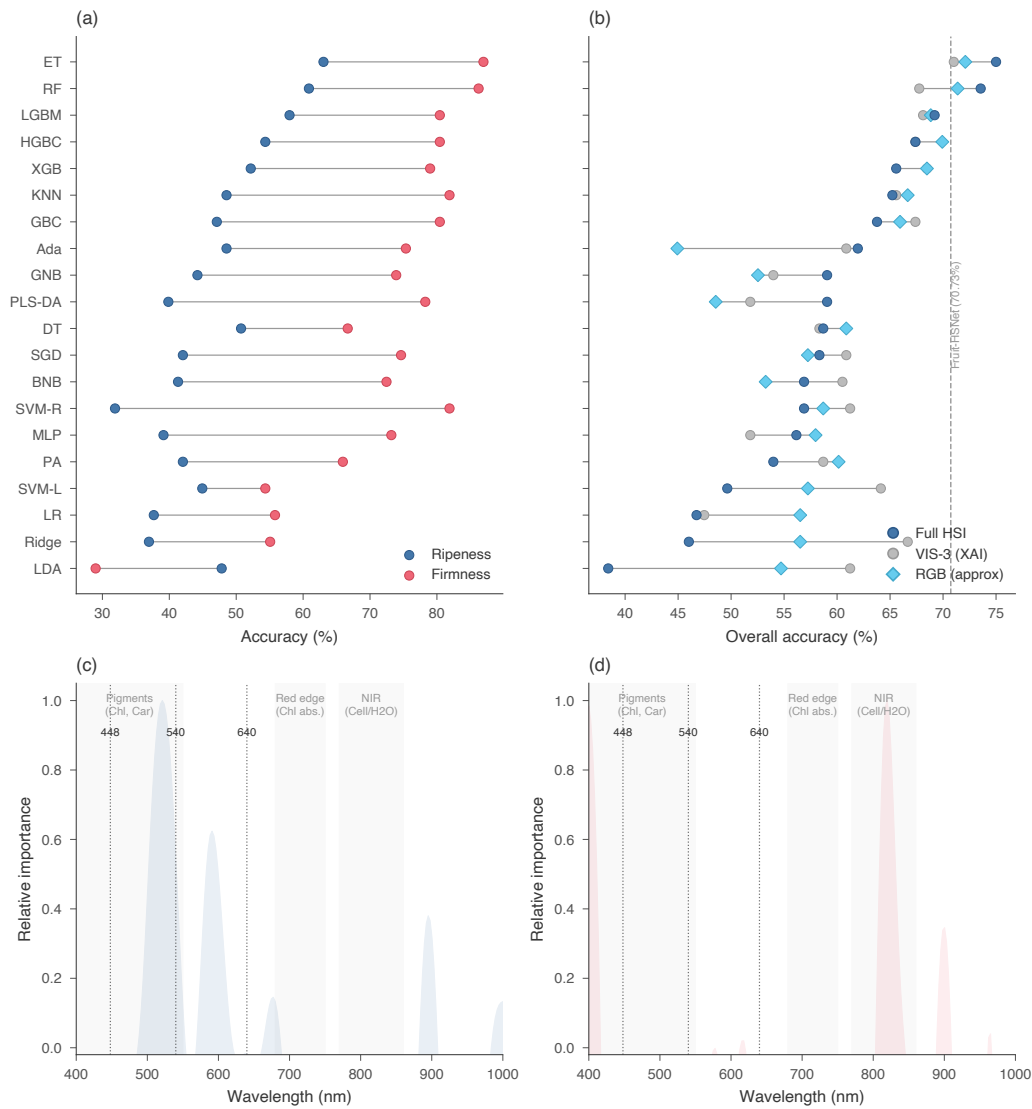


Figure 6: Task asymmetry, modality comparison, and spectral importance for Phase 3 optimized models. **(a)** Task asymmetry: ripeness (pink) vs firmness (green) accuracy per model, revealing a persistent gap across all 20 algorithms. **(b)** Modality comparison: full-spectrum HSI (blue) vs VIS-3 subset (orange) overall accuracy, with the Fruit-HSNet benchmark (70.73%) shown as a dashed line. **(c, d)** Spectral band importance (rolling mean  $\pm$  std) for ExtraTrees across 400–1000 nm; shaded regions highlight biologically meaningful zones (pigment changes, red-edge chlorophyll absorption, NIR cellular/water signatures). Ripeness importance concentrates in visible wavelengths (c), while firmness importance shifts toward NIR features (d).

502 4.4. Per-Fruit, Per-Class, and Ensemble Analysis

503 Table 6 reports per-fruit accuracy for the top five models (Fig. A.3 in Ap-  
 504 pendix), revealing substantial variation across species. Avocado and mango  
 505 achieved highest accuracies (up to 75.0% for mango with RandomForest),  
 506 while papaya consistently underperformed (33.3–44.4%), likely due to its  
 507 smallest sample size ( $n=51$ ).

Table 6: Per-fruit ripeness accuracy (%) for the top-five models (Phase 4 single-split holdout,  $n=138$ ; ranked by Phase 3 overall accuracy). **Bold** = best per fruit; underline = second-best.

Fruit ( $n$ )	ExtraTrees	RandomForest	LightGBM	HistGradBoosting	XGBoost
Avocado (170)	<b>67.9</b>	<u>65.4</u>	44.4	48.1	38.3
Kaki (68)	<b>41.7</b>	<b>41.7</b>	33.3	<b>41.7</b>	<b>41.7</b>
Kiwi (162)	58.3	54.2	<u>62.5</u>	<b>70.8</b>	50.0
Mango (68)	<u>66.7</u>	<b>75.0</b>	41.7	<u>66.7</u>	50.0
Papaya (51)	<u>33.3</u>	<u>33.3</u>	<u>33.3</u>	<u>33.3</u>	<b>44.4</b>

508 Table 7 provides per-class precision, recall, and F1-score from the single-  
 509 split holdout evaluation. Ripeness classification showed primary confusion  
 510 between adjacent states (unripe/perfect and perfect/overripe), as visible in  
 511 the normalized confusion matrix (Fig. 7a). Stratifying ripeness accuracy by  
 512 firmness group reveals that extreme-firmness samples (soft: 74.1%, very firm:  
 513 73.7%) are classified more accurately than intermediate groups (medium:  
 514 54.4%, firm: 56.5%; Table 7, Fig. 7b), suggesting that fruits at firmness  
 515 extremes present clearer ripeness signatures.

516 Ten-fold stratified cross-validation (Fig. 8, Fig. A.2 in Appendix) con-

517 firmed that top models generalize reliably. ExtraTrees achieved  $82.2 \pm 4.7\%$   
 518 mean OA, with narrow 95% confidence intervals. Firmness accuracy was con-  
 519 sistently higher than ripeness across all models, with tighter variance. This  
 520 pattern reinforces the task asymmetry observed in Phase 3.

Table 7: ExtraTrees Phase 4 single-split holdout analysis ( $n=138$ ). **Top:** ripeness per-class metrics from the confusion matrix (Fig. 7a). **Bottom:** ripeness accuracy stratified by firmness group (4-bin analysis on the original unbalanced split). Extreme-firmness groups achieve higher ripeness accuracy than intermediate groups. \*Unripe precision (1.000) with 47.1% recall means the model rarely predicts unripe; most unripe fruit are misclassified as perfect (41.2%) or overripe (11.8%).

Task	Class	Precision	Recall	F1
Ripeness	Overripe	0.561	<b>0.696</b>	0.621
Ripeness	Perfect	<u>0.569</u>	<u>0.638</u>	<u>0.602</u>
Ripeness	Unripe*	<b>1.000</b>	0.471	<b>0.640</b>
<i>Firmness group</i>		<i>Ripeness accuracy</i>		
	Soft ( $n=27$ )	0.741		
	Medium ( $n=68$ )	0.544		
	Firm ( $n=23$ )	0.565		
	Very firm ( $n=19$ )	<b>0.737</b>		

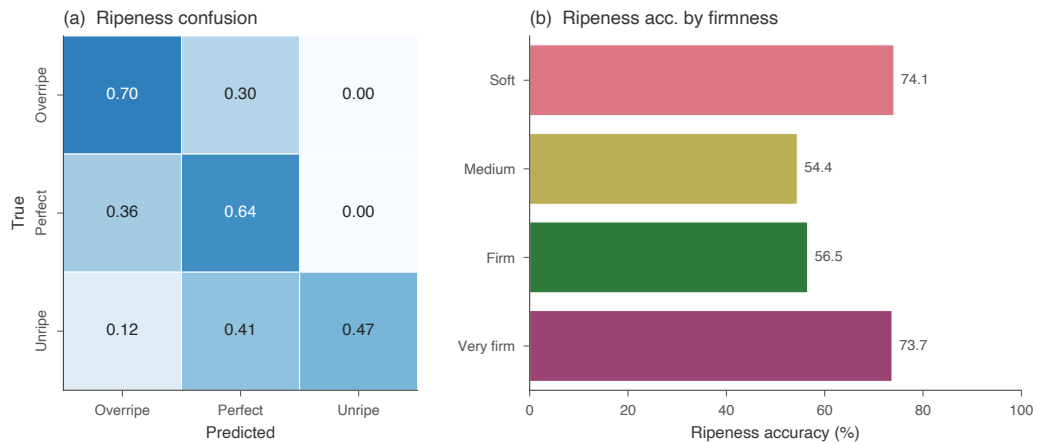


Figure 7: Classification patterns for the best model (ExtraTrees, single-split holdout,  $n=138$ ). **(a)** Normalized ripeness confusion matrix showing primary confusion between adjacent maturation states. **(b)** Ripeness accuracy stratified by firmness group (4-bin); extreme-firmness groups achieve higher ripeness accuracy than intermediate groups.

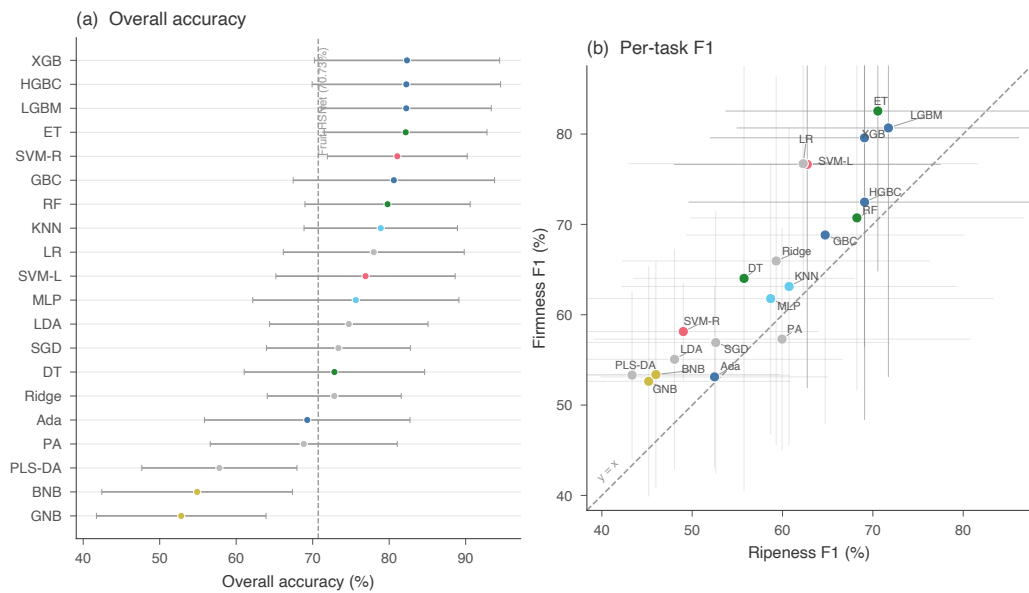


Figure 8: Cross-validation stability (Phase 4, 10-fold stratified CV). **(a)** Overall accuracy with 95% confidence intervals (mean  $\pm t_{9,0.025} \cdot \text{std}$ , where  $t_{9,0.025} \approx 2.262$ ) for all 20 models, sorted by mean. **(b)** Task-separated view: ripeness vs firmness accuracy with CIs. Firmness predictions are consistently more accurate and stable.

521 Ensembles combined the top five base learners using four aggregation  
 522 strategies under original unbalanced and optimal preprocessing configura-  
 523 tions but did not surpass the best single model in either modality (Table 8,  
 524 Fig. 9). The best full-spectrum ensemble (hard/soft voting with stratified re-  
 525 split) achieved 70.65% overall accuracy versus 75.00% for ExtraTrees alone,  
 526 a 4.35 percentage point degradation.

Table 8: Ensemble performance comparison. All ensembles use the top-5 models from Phase 3.  $\Delta$ ET shows difference from best single model (ExtraTrees, 75.0%).

Strategy	Data State	Ripe $\uparrow$	Firm $\uparrow$	OA $\uparrow$	$\Delta$ ET
Hard Voting	Original	52.9	85.5	69.2	-5.8
Soft Voting	Original	50.0	84.8	67.4	-7.6
Stacking	Original	31.2	75.4	53.3	-21.7
Blending	Original	49.3	72.5	60.9	-14.1
Hard Voting	Strat. Resplit	57.2	84.1	70.7	-4.3
Soft Voting	Strat. Resplit	58.0	83.3	70.7	-4.3
Stacking	Strat. Resplit	55.8	81.2	68.5	-6.5
Blending	Strat. Resplit	57.2	82.6	69.9	-5.1

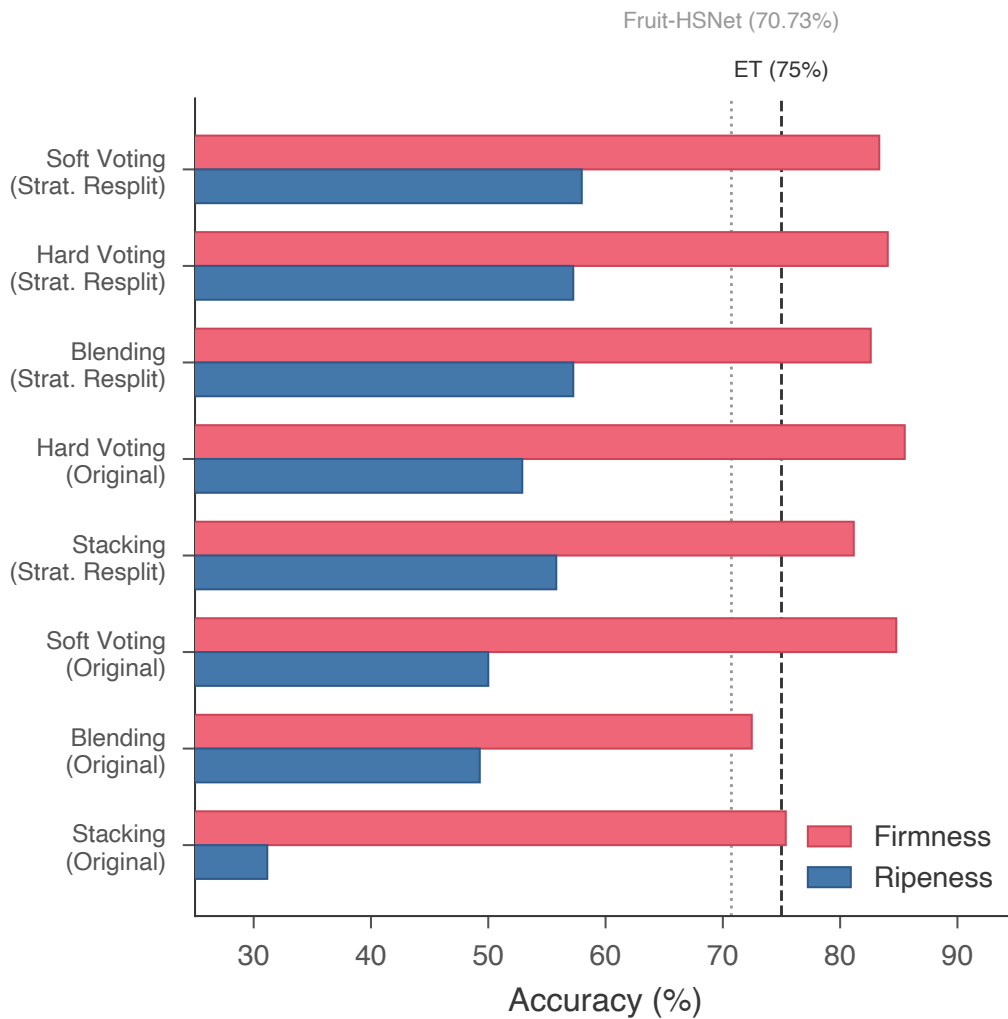


Figure 9: Ensemble performance decomposition by task. All eight ensemble configurations maintain high firmness accuracy ( $>72\%$ ) but catastrophically degrade ripeness accuracy (stacking drops to  $31.2\%$ ). Dashed line: ExtraTrees single-model baseline ( $75.0\%$ ); dotted line: Fruit-HSNet benchmark ( $70.73\%$ ).

527 *4.5. Explainable AI Analysis*

528 Phase 6 XAI analysis for ExtraTrees revealed concentrated spectral band  
529 importance in three regions (Fig. 6c,d): visible blue/green (400–550 nm)  
530 associated with pigment changes, red-edge transitions ( $\sim$ 700–750 nm) indi-  
531 cating vegetation stress, and near-infrared ( $\sim$ 770–860 nm) correlated with  
532 cellular structure. Feature group ablation (Fig. 10) confirmed that the five  
533 spectral transformations contribute differently across tasks and models. For  
534 ripeness, first derivatives are the most critical group for LightGBM (8.3%  
535 drop when removed), while other preprocessing groups are redundant or  
536 harmful for ExtraTrees (removing continuum-removed features improves per-  
537 formance by 10.8%). For firmness, continuum-removed features dominate  
538 ExtraTrees (16.6% drop), while all five groups contribute more uniformly to  
539 LightGBM (largest drop 5.0%). These task- and model-dependent feature  
540 contributions explain why the full 1,120-dimensional feature space outper-  
541 forms PCA reduction: each transformation captures complementary discrim-  
542 inative information that varies by task.

543 Ripeness prediction drew primarily on visible-spectrum features associ-  
544 ated with chlorophyll degradation and carotenoid accumulation, while firm-  
545 ness prediction relied more heavily on near-infrared features correlated with  
546 cellular structure and water content. This wavelength-task correspondence  
547 aligns with established plant physiology knowledge (Khodabakhshian and  
548 Emadi, 2017; Feng et al., 2023). However, XAI attributions may partly re-  
549 flect species-level spectral differences rather than within-species quality-state  
550 variation. The PCA and t-SNE projections (Fig. 1c,d) show that fruit type  
551 dominates the variance structure, so wavelength regions that separate species

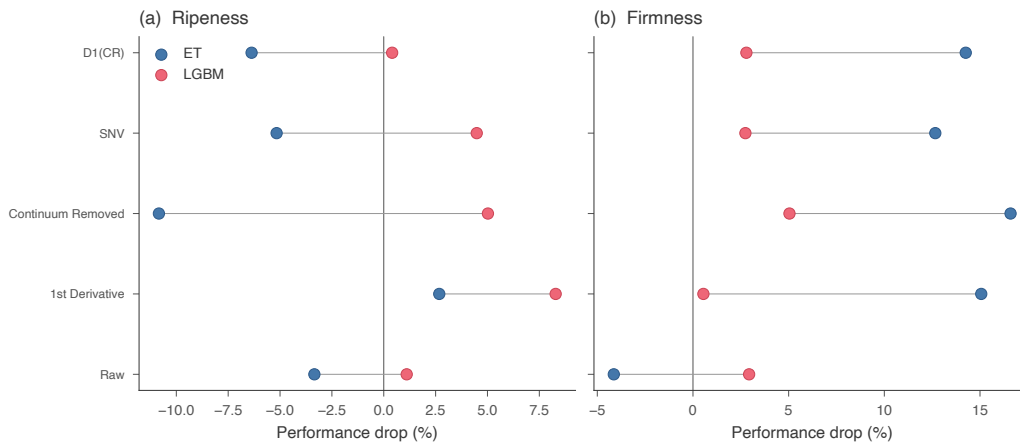


Figure 10: Feature group ablation: relative performance drop (%) when each spectral transformation group (224 bands) is removed. **(a)** Ripeness: first derivatives are most critical for LightGBM (8.3% drop); other groups are redundant or harmful for Extra-Trees. **(b)** Firmness: continuum-removed features dominate ExtraTrees (16.6% drop), while LightGBM shows uniform sensitivity across groups ( $\leq 5.0\%$ ). Positive values indicate removing the group hurts performance; negative values indicate removal improves it. Task- and model-dependent contributions explain why PCA, which compresses all groups uniformly, degrades classification.

552 could receive high importance scores even if they carry limited quality infor-  
553 mation within a single species.

## 554 **5. Discussion**

555 The comprehensive six-phase evaluation reveals three recurring patterns  
556 across 20 algorithms applied to ripeness classification and firmness predic-  
557 tion across five fruit species. Nevertheless, the results raise three fundamental  
558 questions that span all evaluated approaches. What mechanism makes pre-  
559 processing as influential as algorithm choice? Why does the firmness-ripeness  
560 gap persist regardless of model family or modality? And why does PCA, a  
561 standard step in hyperspectral workflows, consistently degrade accuracy on  
562 this feature space? We discuss each in turn before addressing methodological  
563 observations and limitations.

### 564 *5.1. The Interplay of Preprocessing and Algorithm Choice*

565 The relative importance of preprocessing versus algorithm choice depends  
566 on the model family. A bootstrap analysis (10,000 resamples) of the  $20 \times 10$   
567 Phase 2 accuracy matrix yields a 95% CI of  $[-5.8, 21.9]$  pp on the difference  
568 between the two median ranges (algorithm: 32.4 pp; preprocessing: 25.2 pp),  
569 an interval that contains zero. Neither source of variance can be declared the  
570 unambiguous dominant factor with confidence.

571 For the best-performing models, however, preprocessing is the dominant  
572 lever: the best preprocessing configuration improved 15 of 20 models over  
573 their Phase 1 baselines, while Phase 3 optimization did not improve any  
574 model beyond its best Phase 2 result. For weaker models (LDA, Ridge),  
575 preprocessing gains are larger in absolute terms (up to 23 pp) but algorithm

576 choice matters more; these models remain far below competitive accuracy  
 577 regardless of preprocessing. Fig. 11 decomposes each model’s accuracy tra-  
 578 jectory into preprocessing gain (Phase 1→Phase 2) and optimization gain  
 579 (Phase 2→Phase 3), directly visualizing this asymmetry.

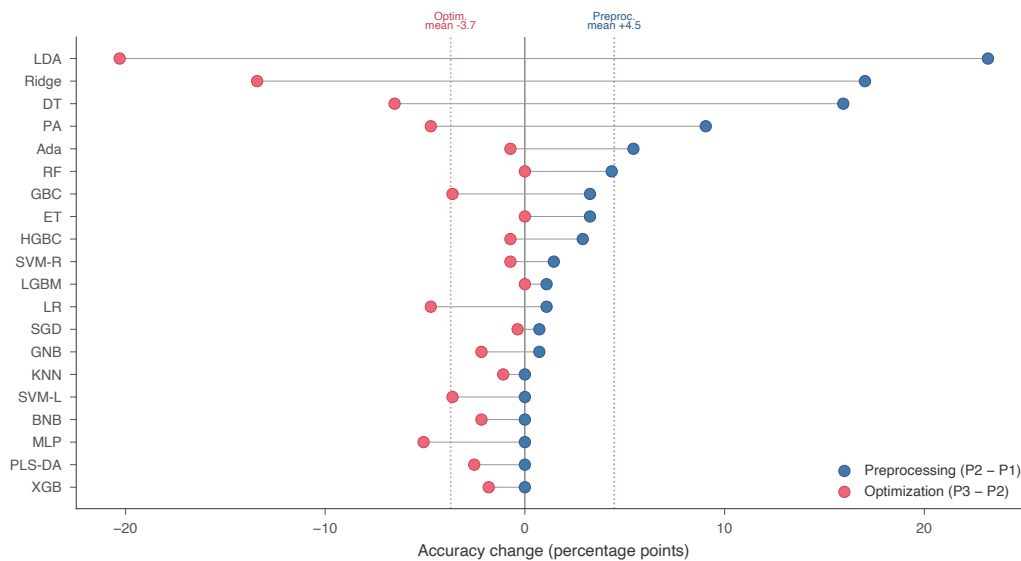


Figure 11: Preprocessing gain vs optimization gain for all 20 models. Each model’s accuracy improvement is decomposed into preprocessing gain (best Phase 2 config minus Phase 1 baseline, teal) and optimization gain (Phase 3 minus best Phase 2, purple). Models sorted by preprocessing gain. Dotted lines indicate means. For most models, preprocessing contributes the dominant share of improvement; optimization gain is near-zero or negative for many.

580 The mechanism behind stratified resplit’s advantage over synthetic re-  
 581 sampling (SMOTE, random oversampling) likely stems from the structure  
 582 of spectral data. SMOTE generates new samples by interpolating between  
 583 existing ones in 1,120-dimensional feature space, an operation that assumes  
 584 local linearity in spectral manifolds shaped by nonlinear biochemical pro-

585 cesses. Interpolating between two spectra at different maturation stages does  
586 not necessarily produce a spectrum corresponding to an intermediate state.  
587 Stratified resplit avoids this problem by redistributing genuine samples across  
588 train and test splits without fabrication. However, the stratified resplit con-  
589 dition also increases training set size by 8.7% (414 vs 381 samples) and im-  
590 proves species balance, so its advantage may partly reflect additional data  
591 and improved species representation rather than purely a balancing effect.

592 The practical implication is that for competitive models, preprocessing in-  
593 vestment yields larger returns than hyperparameter tuning, while algorithm  
594 selection establishes the performance ceiling. Most published studies on this  
595 dataset evaluate a single preprocessing pipeline without systematic ablation.  
596 Our Phase 2 results across 200 configurations (Fig. 5) demonstrate that vary-  
597 ing the preprocessing pipeline alone shifts accuracy by up to 23 pp for the  
598 same algorithm (LDA: 35.5% to 58.7%), which may obscure the true source  
599 of performance variation in published comparisons.

## 600 *5.2. Task Asymmetry as a Construct Validity Question*

601 The persistent firmness–ripeness gap is not simply a matter of task diffi-  
602 culty; it reflects a difference in what the two labels measure.

603 Firmness is a continuous mechanical property with direct spectral cor-  
604 relates. Within the limited NIR range of the FX10 sensor (770–960 nm),  
605 absorption features respond to changes in cellular structure and water con-  
606 tent associated with tissue softening (Lu et al., 2020; Khodabakhshian and  
607 Emadi, 2017). Discretizing firmness into three classes (soft, medium, firm)  
608 creates boundaries that align with real physical thresholds. Ripeness, by con-  
609 trast, is a composite construct. A single ripeness label aggregates chlorophyll

610 degradation, carotenoid accumulation, sugar-acid balance, volatile produc-  
611 tion, and textural changes, processes that occur at different rates and may  
612 produce conflicting spectral signatures within the same fruit. The confusion  
613 between adjacent ripeness states (overripe/perfect: 33.3% mean bidirectional  
614 misclassification, unripe→perfect: 41.2%; Fig. 7a) reflects the inherent am-  
615 biguity of discretizing a multi-dimensional process into ordinal categories.

616 The persistence of this gap across all 20 models and both modalities  
617 suggests the asymmetry is a property of the labeling scheme rather than  
618 of specific algorithms or sensors. Feng et al. (2023) predicted multiple lo-  
619 quat quality parameters and found that prediction accuracy varied across  
620 attributes (color  $R_p^2 = 0.96$ , firmness  $R_p^2 = 0.87$ , SSC  $R_p^2 = 0.84$ ), consis-  
621 tent with the view that different quality dimensions have inherently different  
622 spectral predictability. Future work should investigate whether continuous  
623 regression targets or multi-dimensional maturity indices reduce the informa-  
624 tion loss inherent in categorical ripeness labels.

625 For deployment, the asymmetry is more acceptable than it first appears.  
626 Firmness directly determines shelf life and consumer satisfaction, and com-  
627 mercial distributors prioritize rejecting unacceptably soft fruit over classify-  
628 ing exact ripeness state. However, per-species firmness accuracy (Table 6)  
629 remains below 70% for most individual fruit types at Phase 4, which is in-  
630 sufficient for fully autonomous commercial sorting without human oversight.  
631 The results establish a useful benchmark but should not be interpreted as a  
632 deployable system specification.

### 633 5.3. PCA Degradation and the Variance–Discrimination Mismatch

634 The consistent PCA degradation contradicts standard practice in hyper-  
635 spectral analysis (Lu et al., 2020; Ahmad et al., 2022). The explanation is  
636 rooted in our feature engineering pipeline.

637 The 1,120-feature space concatenates five spectral transformations, each  
638 designed to isolate different physical phenomena: raw reflectance captures  
639 absolute intensity, first derivatives capture absorption edge slopes, contin-  
640 uum removal normalizes band depth, and SNV corrects for scattering. PCA  
641 maximizes variance across the concatenated space, but variance and task-  
642 relevant discrimination are not aligned. First derivatives contribute  $<0.1\%$   
643 of total variance across the 1,120 features, yet removing them drops Light-  
644 GBM ripeness accuracy by 8.3% (Fig. 10a). PCA at 95% retained variance  
645 effectively discards these low-variance, high-discrimination features.

646 The 7 models that selected PCA in their optimal Phase 2 configuration  
647 were predominantly linear (Ridge, LDA, SGDClassifier, PassiveAggressive).  
648 These methods benefit from PCA because dimensionality reduction mitigates  
649 the multicollinearity that destabilizes their coefficient estimates. Tree-based  
650 methods face no such constraint: their split-based decision boundaries are  
651 invariant to feature correlation. This interaction between model family and  
652 preprocessing underlines why blanket preprocessing recommendations of al-  
653 ways applying PCA to hyperspectral data can be counterproductive when  
654 engineered features precede dimensionality reduction.

### 655 5.4. Full-Spectrum Versus Reduced-Band Sensing

656 The full-spectrum advantage at the top of the leaderboard is modest:  
657 ExtraTrees achieves 75.0% with 1,120 features compared to 72.1% with ap-

658 proximate RGB bands (96.1% recovery) and 71.0% with XAI-derived VIS-3  
659 bands (94.7% recovery). Hyperspectral systems cost \$10,000–100,000 com-  
660 pared to \$500–5,000 for multispectral or RGB alternatives (Min et al., 2023).  
661 Both reduced-band configurations recover 94–96% of full-spectrum accuracy  
662 using 15 features, and several models achieve comparable or better Phase 3  
663 OA in the low-dimensional space (Table 5).

664 The VIS-3 and RGB configurations produce statistically indistinguish-  
665 able results across 20 models (mean OA 61.7% vs 60.1%; Cohen’s  $d = 0.27$ ,  
666  $p = 0.30$ ; Table 3). The two band sets differ by only 1–4 spectral indices  
667 (2–11 nm in wavelength), yet neither configuration holds a consistent advan-  
668 tage. This insensitivity to exact band placement indicates that the discrim-  
669 inative spectral content in the visible range is spectrally broad: any three  
670 well-separated bands within the 400–700 nm window capture it. XAI-guided  
671 band selection provides no measurable benefit over generic RGB-center wave-  
672 lengths on this benchmark.

673 XAI analysis concentrates discriminative wavelengths in three regions:  
674 visible blue/green (400–550 nm) for pigment assessment, the red edge ( $\sim 700$ –  
675 750 nm) for chlorophyll transitions, and near-infrared ( $\sim 770$ –860 nm) for  
676 cellular structure (Fig. 6c,d). This concentration aligns with the spectral-  
677 biochemical relationships established by Khodabakhshian and Emadi (2017)  
678 and Feng et al. (2023). However, interpretations specific to the short-NIR  
679 boundary (770–960 nm) should be treated as speculative relative to the  
680 well-established core NIR firmness bands (1200–1450 nm), which fall out-  
681 side the Specim FX10 range; the firmness-related importance observed here  
682 may partly capture species-level variance rather than within-species quality

683 differences.

684 Two caveats apply. First, both reduced-band configurations extract fea-  
685 tures from hyperspectral cubes at exact wavelengths, bypassing the spectral  
686 response curves, Bayer filter artifacts, and noise characteristics of real sen-  
687 sors. The reduced-band accuracy figures represent upper bounds. That said,  
688 the RGB parity result weakens the “upper bound” framing: if approximate  
689 RGB wavelengths match XAI-optimized bands, the gap between simulated  
690 and real sensor performance may be smaller than the VIS-3 result alone sug-  
691 gests. Second, the full-spectrum advantage may increase for tasks requiring  
692 finer spectral discrimination, such as early-stage decay detection where subtle  
693 NIR absorption shifts precede visible symptoms (Min et al., 2023).

#### 694 5.5. *Tree-Based Dominance and Ensemble Underperformance*

695 Tree-based ensembles (ExtraTrees, RandomForest) occupied the top two  
696 positions in Phase 3 and remained competitive in Phase 4 cross-validation  
697 ( $82.2 \pm 4.7\%$  and  $79.8 \pm 4.8\%$ ). Their advantage arises from a structural  
698 match with the feature space. The 1,120 features comprise five transforma-  
699 tion groups (224 bands each) with high within-group correlation but distinct  
700 cross-group discriminative content (Fig. 10). ExtraTrees’ random split-point  
701 selection (Geurts et al., 2006) acts as implicit random subspace sampling  
702 across these groups, decorrelating individual trees without explicit feature  
703 selection. Linear methods face a different problem: multicollinearity among  
704 the 224 correlated bands within each transformation inflates coefficient vari-  
705 ance and degrades generalization.

706 The ensemble underperformance with homogeneous base learners (Ta-  
707 ble 8) is instructive. All four aggregation strategies degraded performance

708 relative to ExtraTrees alone, with stacking suffering the worst drop ( $-21.7$  pp  
709 under original data). The top-five base learners (ExtraTrees, RandomForest,  
710 LGBM, HistGradientBoosting, XGBoost) are all recursive partitioning meth-  
711 ods that construct decision boundaries through axis-aligned splits. They  
712 make correlated errors, particularly on the same borderline ripeness cases  
713 (Table 7). Ensemble theory predicts that aggregation improves predictions  
714 only when base learner errors are sufficiently uncorrelated. Including a struc-  
715 turally different learner (e.g., SVM or a neural network) as a base learner  
716 might restore the diversity that ensemble methods require, though at the  
717 cost of interpretability.

718 ExtraTrees, despite providing feature importance scores through tree-  
719 based splitting, is not fully transparent: predictions for individual samples  
720 cannot be explained without post-hoc XAI tools. The same limitation ap-  
721 plies to all gradient-boosted methods evaluated here. The interpretability  
722 advantage of tree-based methods is relative to black-box neural networks,  
723 not absolute.

#### 724 *5.6. Spatial Complexity and Mean-Spectrum Representations*

725 ExtraTrees achieves 75.0% on the benchmark split using only per-sample  
726 mean spectra, a representation that retains spectral wavelength information  
727 but discards all spatial variation across fruit pixels. Fruit-HSNet’s dual-  
728 branch architecture (Section 3) processes full spatial cubes yet reports 70.73%  
729 under a different protocol. Two interpretations are consistent with this out-  
730 come. Spatial variation across fruit pixels may carry little additional discrim-  
731 inative signal beyond what the mean spectrum encodes. Alternatively, mean  
732 spectrum aggregation may lose useful spatial information, but any advantage

733 Fruit-HSNet derives from recovering it may be offset by other architectural  
734 differences. Distinguishing these interpretations requires ablation within the  
735 Fruit-HSNet framework, which the available data cannot support. The prac-  
736 tical implication holds regardless of cause. ExtraTrees trains in 0.2 seconds  
737 on a consumer CPU, requires no data augmentation, produces a 2.77 MB  
738 model, and achieves 2.1 ms per-sample inference latency.

739 The three discriminative wavelengths identified by explainability analysis  
740 (448, 540, 640 nm) were selected as the top joint-ranked visible-range bands  
741 from ExtraTrees consensus feature importance across both tasks (Fig. 6c,d,  
742 dotted lines). These wavelengths lie within 2–11 nm of standard RGB filter  
743 centers (450, 550, 650 nm), and the controlled comparison confirms that this  
744 proximity translates to equivalent classification performance: RGB-center  
745 bands recover 96.1% of full-spectrum accuracy versus 94.7% for VIS-3 (Ta-  
746 ble 5). Commodity RGB cameras use broad, overlapping passbands rather  
747 than narrow spectral selections, but the insensitivity to exact band place-  
748 ment observed here suggests that such cameras may capture sufficient dis-  
749 criminative content for this task. Validation with actual RGB sensor data  
750 remains necessary to confirm this implication, because Bayer-pattern demo-  
751 saicing, sensor noise, and spectral crosstalk introduce artifacts absent from  
752 simulated band selection.

753 The spectroscopic basis of these feature importance findings carries direct  
754 implications for purpose-built inline sensing hardware. Discriminative infor-  
755 mation concentrated near the chlorophyll absorption trough ( $\sim 680$  nm), the  
756 red-edge transition (700–750 nm), and the tissue water/structure region (770–  
757 860 nm) suggests that a targeted multispectral filter array covering these

758 regions could capture a substantial portion of the discriminative spectral  
759 content available in the full 224-band dataset, at a fraction of hyperspectral  
760 system cost. Reducing to 6–12 strategically positioned channels would lower  
761 per-sample data volume sufficiently to enable inference on embedded con-  
762 trollers without GPU acceleration, a key deployment constraint in high-speed  
763 sorting environments. The demonstrated importance of SNV normalization,  
764 illustrated by the 12.7% relative performance drop in ExtraTrees firmness  
765 upon its removal (Fig. 10b), is not merely a preprocessing detail but a spec-  
766 troscopically motivated robustness measure against the multiplicative scatter  
767 effects of variable surface geometry and path length inherent to production-  
768 line acquisition. On a reduced-band sensor, however, computing SNV statis-  
769 tics across only a handful of spectral points becomes poorly conditioned; al-  
770 ternative scatter correction strategies such as external diffuse reference tiles  
771 or integrating sphere-based calibration panels would be required to preserve  
772 this robustness in practice.

773 Push-broom line-scan sensors are architecturally well-matched to con-  
774 veyor belt translation, but rotating or irregularly presented fruit introduces  
775 variable illumination geometry that complicates line-scan acquisition. Snap-  
776 shot multispectral mosaic sensors, which acquire full spatial and spectral  
777 information in a single frame, may be better suited to this geometry at the  
778 cost of lower spectral resolution, and represent a complementary hardware  
779 direction worth evaluating against the wavelength importance findings re-  
780 ported here.

781 *5.7. Single-Split Versus Cross-Validated Evaluation*

782 A pattern in Table 1 warrants attention: several models rank very dif-  
783 ferently in Phase 3 (single train-test split) versus Phase 4 (10-fold CV).  
784 SVC\_RBF achieves only 56.9% in Phase 3 but  $81.1 \pm 4.0\%$  in Phase 4. Logis-  
785 ticRegression drops to 46.7% in Phase 3 while reaching  $78.0 \pm 5.2\%$  in Phase 4.  
786 Conversely, GaussianNB scores 59.1% in Phase 3 but only  $52.8 \pm 4.9\%$  in  
787 Phase 4.

788 These discrepancies arise because Phase 3 evaluates on a single held-out  
789 set, making results sensitive to the specific composition of that split. Phase 4  
790 cross-validation averages over 10 splits and provides confidence intervals. The  
791 models most affected (SVMs, LogisticRegression, LDA) are those whose deci-  
792 sion boundaries are sensitive to the training sample composition. Tree-based  
793 ensembles show smaller discrepancies because bagging and random subspace  
794 sampling provide internal averaging that reduces split sensitivity. Even with  
795 identical hyperparameters and data splits, stochastic models such as Ex-  
796 traTrees may vary across re-trainings due to internal randomization (e.g.,  
797 random subspace selection); the  $\sim 2$  pp gap between ExtraTrees’ Phase 3  
798 holdout accuracy (75.0%) and Phase 4 holdout re-evaluation (73.2%) reflects  
799 this inherent variance rather than a methodological inconsistency.

800 This observation has a methodological implication for the field: studies  
801 reporting accuracy from a single train-test partition may rank algorithms  
802 differently than cross-validated evaluations. The Fruit-HSNet benchmark  
803 (70.73%) was itself evaluated on a single split (Ben Jmaa et al., 2025), which  
804 limits the precision of any direct comparison. Our Phase 4 cross-validated  
805 estimates are more reliable for algorithm ranking, though they evaluate a

806 different quantity (generalization across random partitions) than Phase 3  
807 (performance on the dataset’s original held-out set).

### 808 *5.8. Limitations*

809 Four limitations affect the generalizability of these findings most directly.  
810 First, all results derive from a single dataset (DeepHS Fruit) captured with  
811 one camera system (Specim FX10) under laboratory conditions across 47  
812 acquisition days. Field environments introduce illumination variation, back-  
813 ground clutter, and 15–30% accuracy degradation (Min et al., 2023). Whether  
814 the preprocessing-over-algorithm ordering holds under field conditions re-  
815 mains untested. Per-fruit accuracy (33–67% for individual species at Phase 4)  
816 falls well below the threshold required for fully autonomous commercial sort-  
817 ing without human oversight; these results should be interpreted as a research  
818 benchmark, not a deployment claim.

819 Second, the comparison to Fruit-HSNet (70.73%) involves a protocol mis-  
820 match: Ben Jmaa et al. (2025) used a different train-test partition and eval-  
821 uation procedure. The Phase 3 result (ExtraTrees, 75.0%; Wilson 95% CI:  
822 [67.6%, 81.8%]) and Fruit-HSNet’s reported value are not directly compa-  
823 rable. The ExtraTrees confidence interval contains the Fruit-HSNet value.  
824 The scikit-learn `MLPClassifier` benchmarked here is a shallow feedforward  
825 network, not a deep learning model. ExtraTrees’ 75.0% numerically exceeds  
826 Fruit-HSNet’s 70.73%, but the CI contains that value and protocol differ-  
827 ences preclude a controlled comparison; no general superiority claim over  
828 deep learning architectures is warranted. The benchmark is scoped to the  
829 scikit-learn ecosystem.

830 Third, sample sizes are imbalanced across fruit species (papaya  $n=51$  vs

831 avocado  $n=170$ ), and papaya consistently underperforms (33–44% accuracy;  
832 Table 6). Per-fruit deployment would require species-specific calibration and  
833 substantially larger training sets for underrepresented varieties.

834 Fourth, PLS-DA, the standard chemometric baseline, is included but  
835 ranks last among the 20 models (P4 CV:  $57.8 \pm 4.5\%$ ), likely because the  
836 effective number of latent components is clamped to the number of classes  
837 ( $n = 3$ ), severely limiting its capacity on the 1,120-feature space. This re-  
838 sult should not be interpreted as evidence that PLS-DA is inherently inferior  
839 for spectral classification; with larger class counts or continuous regression  
840 targets, PLS-DA may perform competitively.

## 841 **6. Conclusion**

842 In this study, we investigated various preprocessing strategies and ma-  
843 chine learning algorithms for non-destructive estimation of fruit ripeness and  
844 firmness from hyperspectral imaging. Our evaluations showed that prepro-  
845 cessing choices can influence hyperspectral fruit quality prediction as strongly  
846 as algorithm selection, with preprocessing improvements sometimes exceed-  
847 ing gains from model optimization. The consistent performance gap between  
848 firmness and ripeness prediction appears to arise from differences in label  
849 structure rather than algorithmic limitations. Tree-based ensemble models  
850 achieved up to 82.2% cross-validated accuracy on standard consumer CPU  
851 hardware. Additionally, explainable AI analysis identified three informative  
852 visible-range wavelengths that retained 94.7% of full-spectrum performance  
853 using only 15 features, highlighting their potential for targeted multispectral  
854 sensor design. Future work should validate these findings across larger multi-

855 species datasets to assess the generalizability of the preprocessing–algorithm  
 856 relationship and the transferability of the identified wavelengths.

857 **Appendix A. Supplementary Figures**

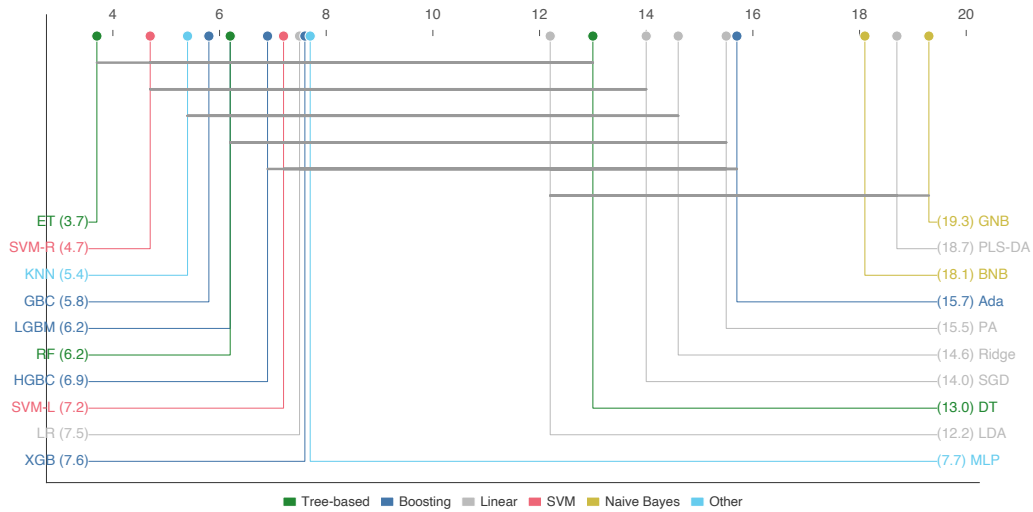


Figure A.1: Critical difference diagram from Nemenyi post-hoc test ( $\alpha = 0.05$ ) on 20-model  $\times$  10-fold overall F1 scores. Models connected by a horizontal bar are not significantly different. Models colored by algorithm family.

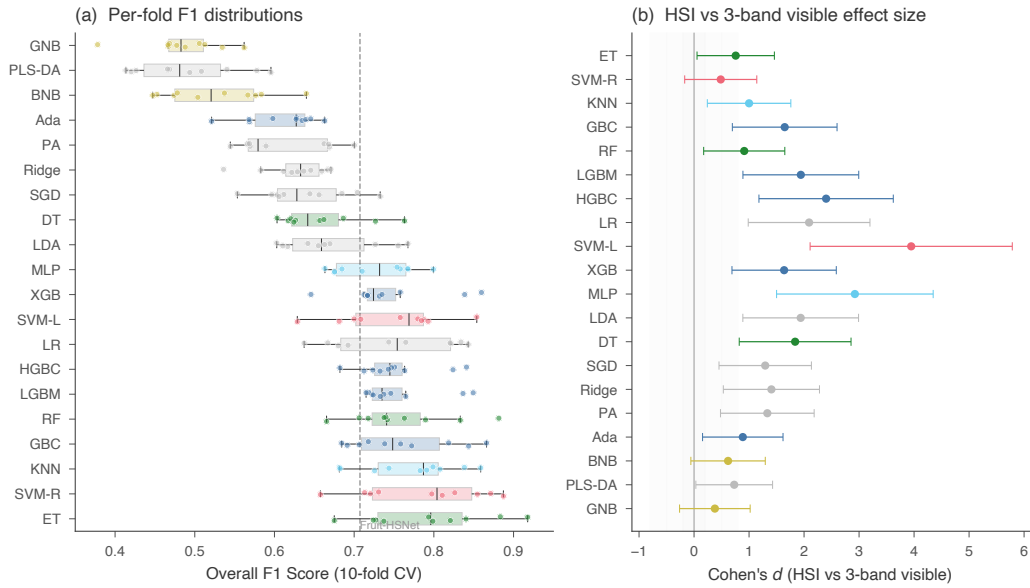


Figure A.2: (a) Per-fold overall F1 score distributions from 10-fold stratified cross-validation, sorted by mean F1. Box plots show median and interquartile range; individual fold scores are overlaid as jittered points. Dashed line marks the Fruit-HSNet benchmark (70.7%). (b) Per-model Cohen's  $d$  effect sizes for the HSI vs. VIS-3 subset modality comparison, with 95% confidence intervals. Shaded bands indicate negligible ( $|d| < 0.2$ ), small ( $0.2-0.5$ ), and medium ( $0.5-0.8$ ) effect regions. Models colored by algorithm family.

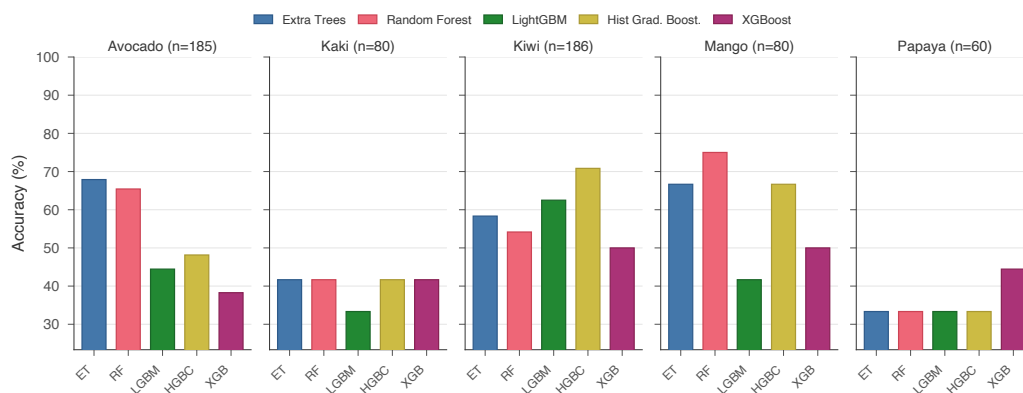


Figure A.3: Per-fruit overall accuracy for the top-five models (Phase 4). Avocado and mango show highest performance; papaya consistently underperforms due to smallest sample size ( $n=51$ ). Performance variation across species highlights deployment considerations for multi-fruit sorting systems.

858 **Appendix B. Supplementary Tables**

Table B.1: Overview of 20 evaluated models with Optuna hyperparameter search spaces. HP = number of tuned hyperparameters; FI = native feature importance. All ranges define the Bayesian optimization bounds (100 trials per model).

Abbr.	Model	HP	FI	Lib.	Optuna Search Space
ET	Extra Trees	5	✓	sklearn	$n\_est \in [50, 500]$ , $depth \in [3, 30]$ , $split \in [2, 20]$ , $leaf \in [1, 20]$ , $crit \in \{gini, ent\}$
RF	Random Forest	6	✓	sklearn	$n\_est \in [50, 500]$ , $depth \in [3, 30]$ , $split \in [2, 20]$ , $leaf \in [1, 20]$ , $feat \in \{\text{sqrt}, \log 2, \emptyset\}$ , $crit \in \{gini, ent\}$
DT	Decision Tree	4	✓	sklearn	$depth \in [3, 30]$ , $split \in [2, 20]$ , $leaf \in [1, 20]$ , $crit \in \{gini, ent\}$
XGB	XGBoost	9	✓	xgboost	$n\_est \in [50, 500]$ , $lr \in [.01, .3]$ , $depth \in [3, 10]$ , $child\_w \in [1, 10]$ , $\gamma \in [0, 1]$ , $sub \in [.5, 1]$ , $col \in [.5, 1]$ , $\alpha \in [0, 1]$ , $\lambda \in [0, 1]$
LGBM	LightGBM	8	✓	lightgbm	$n\_est \in [50, 500]$ , $lr \in [.01, .3]$ , $depth \in [3, 15]$ , $leaves \in [10, 300]$ , $child \in [5, 100]$ , $\alpha \in [0, 1]$ , $\lambda \in [0, 1]$ , $col \in [.5, 1]$
GBC	Gradient Boosting	5	✓	sklearn	$n\_est \in [50, 500]$ , $lr \in [.01, .3]$ , $depth \in [3, 10]$ , $split \in [2, 20]$ , $leaf \in [1, 20]$
HGBC	HistGradBoosting	5	-	sklearn	$lr \in [.01, .3]$ , $iter \in [100, 500]$ , $depth \in [3, 10]$ , $leaf \in [10, 100]$ , $L2 \in [0, 1]$
Ada	AdaBoost	3	✓	sklearn	$n\_est \in [25, 200]$ , $lr \in [.1, 1.5]$ , $alg = \text{SAMME}$
LR	Logistic Regression	3	-	sklearn	$C \in [.01, 100]$ , $pen \in \{L1, L2\}$ , $solver \in \{\text{liblinear}, \text{saga}\}$
Ridge	Ridge Classifier	2	-	sklearn	$\alpha \in [.1, 100]$ , $solver \in \{\text{auto}, \text{svd}, \text{cholesky}, \text{lsqr}, \text{sag}, \text{saga}\}$
SGD	SGD Classifier	4	-	sklearn	$\alpha \in [1e-6, .01]$ , $pen \in \{L1, L2\}$ , $lr \in \{\text{const}, \text{opt}, \text{inv}, \text{adapt}\}$ , $\eta_0 \in [.01, 1]$
PA	Passive-Aggressive	3	-	sklearn	$C \in [.01, 10]$ , $iter \in [100, 2k]$ , $tol \in [1e-5, .01]$
LDA	Linear Discriminant	2	-	sklearn	$solver \in \{\text{svd}, \text{lsqr}, \text{eigen}\}$ , $shrink \in \{\emptyset, \text{auto}, .1, .5, .9\}$
SVM-L	SVC (Linear)	3	-	sklearn	$C \in [.01, 100]$ , $tol \in [1e-5, .01]$ , $iter \in [1k, 5k]$
SVM-R	SVC (RBF)	3	-	sklearn	$C \in [.01, 100]$ , $tol \in [1e-5, .01]$ , $\gamma \in [1e-6, .1]$
GNB	Gaussian NB	1	-	sklearn	$var\_smooth \in [1e-10, 1e-6]$
BNB	Bernoulli NB	2	-	sklearn	$\alpha \in [.1, 10]$ , $fit\_prior \in \{T, F\}$
KNN	k-NN	4	-	sklearn	$k \in [3, 50]$ , $w \in \{\text{uniform}, \text{dist}\}$ , $alg \in \{\text{auto}, \text{ball}, \text{kd}, \text{brute}\}$ , $p \in \{1, 2\}$
MLP	MLP	4	-	sklearn	$hidden \in \{(50), (100), (50, 50), (100, 50), (100, 100)\}$ , $\alpha \in [1e-4, .1]$ , $lr \in [.001, .1]$ , $act \in \{\text{relu}, \text{tanh}, \text{logistic}\}$
PLS-DA	PLS Discriminant Analysis	1	-	sklearn	$n\_comp \in [1, 30]$ (clamped to $\min(n_{\text{samples}}, n_{\text{features}}, n_{\text{classes}})$ )

Table B.2: Best hyperparameters found by Optuna (100 Bayesian TPE trials per model). Parameters listed after stripping the `classifier__` prefix. Models ranked by Phase 3 overall accuracy.

Rk	Model	Best Hyperparameters
1	ExtraTrees	<code>n_estimators=277, max_depth=14, min_samples_split=6, min_samples_leaf=1, criterion=gini</code>
2	RandomForest	<code>n_estimators=285, max_depth=10, min_samples_split=5, min_samples_leaf=2, max_features=log2, criterion=entropy</code>
3	LGBM	<code>n_estimators=470, learning_rate=0.2874, max_depth=14, num_leaves=10, min_child_samples=90, reg_alpha=0.2373, reg_lambda=0.4189, colsample_bytree=0.6112</code>
4	HistGradientBoosting	<code>learning_rate=0.1795, max_iter=123, max_depth=5, min_samples_leaf=41, l2_regularization=0.5475</code>
5	XGBoost	<code>n_estimators=400, learning_rate=0.1357, max_depth=10, min_child_weight=8, gamma=0.1722, subsample=0.7181, colsample_bytree=0.5824, reg_alpha=0.3973, reg_lambda=0.5746</code>
6	KNeighbors	<code>n_neighbors=3, weights=distance, algorithm=ball_tree, p=1</code>
7	GradientBoosting	<code>n_estimators=181, learning_rate=0.1839, max_depth=8, min_samples_split=11, min_samples_leaf=17</code>
8	AdaBoost	<code>n_estimators=198, learning_rate=0.9135, algorithm=SAMME</code>
9	GaussianNB	<code>var_smoothing=3.15e-09</code>
10	DecisionTree	<code>max_depth=8, min_samples_split=17, min_samples_leaf=8, criterion=entropy</code>
11	SGDClassifier	<code>penalty=l2, alpha=5.18e-04, learning_rate=adaptive, eta0=0.3370</code>
12	BernoulliNB	<code>alpha=0.1026, fit_prior=True</code>
13	SVC_RBF	<code>C=30.96, tol=1.17e-05, gamma=0.0019</code>
14	MLPClassifier	<code>hidden_layer_sizes=(100, 50), alpha=0.0255, learning_rate_init=0.0017, activation=relu</code>
15	PassiveAggressive	<code>C=3.85, max_iter=369, tol=0.0092</code>
16	SVC_Linear	<code>C=0.0880, tol=1.56e-05, max_iter=2560</code>
17	LogisticRegression	<code>penalty=l2, C=35.64, solver=liblinear</code>
18	Ridge	<code>alpha=93.90, solver=cholesky</code>
19	LDA	<code>solver=lsqr, shrinkage=0.1000</code>
20	PLSDA	<code>n_components=2</code>

Table B.3: Hardware and Software Environment Specifications. All 20 models were trained on CPU only; no GPU acceleration was used.

<b>Component</b>	<b>Specification</b>
<i><b>Hardware</b></i>	
Processor	Apple M4 (10-core CPU: 4 performance + 6 efficiency)
Memory	24 GB unified LPDDR5X
Storage	SSD
GPU	None (CPU-only training)
<i><b>Software Environment</b></i>	
Operating System	macOS (arm64)
Python	3.13
scikit-learn	1.7.1
XGBoost	3.0.4
LightGBM	4.6.0
Optuna	4.5.0
imbalanced-learn	0.14.0
SHAP	0.48.0
LIME	0.2.0.1
<i><b>Reproducibility Settings</b></i>	
Global random seed	42
Seeded components	Python, NumPy, all classifiers, SMOTE, splits, Optuna TPE sampler
Data split	Dataset-provided VIS benchmark split (73.4/26.6)
Cross-validation	StratifiedKFold, $k=10$ , shuffle=True
Optuna trials	100 per model (TPE sampler, seed=42)
Thread pinning	OMP/OpenBLAS/MKL = 1 thread

859 **Data Availability**

860 DeepHS Fruit dataset (version 2) (Varga et al., 2021) is available at  
861 [https://github.com/cogsys-tuebingen/deephs\\_fruit](https://github.com/cogsys-tuebingen/deephs_fruit). Code and sup-  
862plementary materials will be made available upon publication.

863 **Declaration of Competing Interest**

864 The authors declare that they have no known competing financial inter-  
865 ests or personal relationships that could have appeared to influence the work  
866 reported in this paper.

867 **CRedit Authorship Contribution Statement**

868 **Phongsakon Mark Konrad:** Conceptualization, Methodology, Soft-  
869 ware, Validation, Formal analysis, Investigation, Data curation, Writing –  
870 original draft, Visualization. **Casper Kunstmann-Olsen:** Supervision,  
871 Writing – review & editing. **Jacek Fiutowski:** Supervision, Writing – re-  
872 view & editing. **Serkan Ayvaz:** Conceptualization, Supervision, Writing –  
873 review & editing, Project administration.

874 **Acknowledgements**

875 This research did not receive any specific grant from funding agencies in  
876 the public, commercial, or not-for-profit sectors.

877 **Declaration of Generative AI and AI-assisted Technologies in the**  
878 **Writing Process**

879 During the preparation of this work the authors used Claude Opus 4.6 and  
880 Claude Sonnet 4.6 (Anthropic) for writing assistance and code development,  
881 and Claude Haiku 4.5 (Anthropic) for routine code formatting tasks. After  
882 using these tools, the authors reviewed and edited the content as needed and  
883 take full responsibility for the content of the published article.

884 **References**

885 Ahmad, M., Shabbir, S., Roy, S.K., Hong, D., Wu, X., Yao, J., Khan, A.M.,  
886 Mazzara, M., Distefano, S., Chanussot, J., 2022. Hyperspectral image  
887 classification—traditional to deep models: A survey for future prospects.  
888 IEEE Journal of Selected Topics in Applied Earth Observations and Re-  
889 mote Sensing 15, 968–999. doi:10.1109/JSTARS.2021.3133021.

890 Ahmed, M.T., Monjur, O., Kamruzzaman, M., 2024a. Deep learning-based  
891 hyperspectral image reconstruction for quality assessment of agro-product.  
892 Journal of Food Engineering 382, 112223. doi:10.1016/j.jfoodeng.2024.  
893 112223.

894 Ahmed, M.T., Wijewardane, N.K., Lu, Y., Jones, D.S., Kudenov, M.,  
895 Williams, C., Villordon, A., Kamruzzaman, M., 2024b. Advancing sweet-  
896 potato quality assessment with hyperspectral imaging and explainable ar-  
897 tificial intelligence. Computers and Electronics in Agriculture 220, 108855.  
898 doi:10.1016/j.compag.2024.108855.

- 899 Akiba, T., Sano, S., Yanase, T., Ohta, T., Koyama, M., 2019. Optuna: A  
900 next-generation hyperparameter optimization framework, in: Proceedings  
901 of the 25th ACM SIGKDD International Conference on Knowledge Dis-  
902 covery and Data Mining, ACM. pp. 2623–2631. doi:10.1145/3292500.  
903 3330701.
- 904 Ben Jmaa, A.B., Chaieb, F., Fabijańska, A., 2025. Fruit-hsnet: A machine  
905 learning approach for hyperspectral image-based fruit ripeness prediction,  
906 in: Proceedings of the 17th International Conference on Agents and Ar-  
907 tificial Intelligence (ICAART), SCITEPRESS. pp. 102–111. doi:10.5220/  
908 0013110800003890.
- 909 Bergstra, J., Bardenet, R., Bengio, Y., Kégl, B., 2011. Algorithms for hyper-  
910 parameter optimization, in: Advances in Neural Information Processing  
911 Systems, pp. 2546–2554. doi:10.5555/2986459.2986743.
- 912 Breiman, L., 2001. Random forests. *Machine Learning* 45, 5–32. doi:10.  
913 1023/A:1010933404324.
- 914 Chen, T., Guestrin, C., 2016. XGBoost: A scalable tree boosting system,  
915 in: Proceedings of the 22nd ACM SIGKDD International Conference on  
916 Knowledge Discovery and Data Mining, ACM. pp. 785–794. doi:10.1145/  
917 2939672.2939785.
- 918 Feng, S., Shang, J., Tan, T., Wen, Q., Meng, Q., 2023. Nondestructive quality  
919 assessment and maturity classification of loquats based on hyperspectral  
920 imaging. *Scientific Reports* 13, 13189. doi:10.1038/s41598-023-40553-3.

- 921 Gao, Z., Shao, Y., Xuan, G., Wang, Y., Liu, Y., Han, X., 2020. Real-time  
922 hyperspectral imaging for the in-field estimation of strawberry ripeness  
923 with deep learning. *Artificial Intelligence in Agriculture* 4, 31–38. doi:10.  
924 1016/j.aiia.2020.04.003.
- 925 Garillos-Manliguez, C.A., Chiang, J.Y., 2021. Multimodal deep learning  
926 and visible-light and hyperspectral imaging for fruit maturity estimation.  
927 *Sensors* 21, 1288. doi:10.3390/s21041288.
- 928 Geurts, P., Ernst, D., Wehenkel, L., 2006. Extremely randomized trees.  
929 *Machine Learning* 63, 3–42. doi:10.1007/s10994-006-6226-1.
- 930 Guo, Z., Zhang, J., Wang, H., Li, S., Shao, X., Dong, H., Sun, J., Geng,  
931 L., Zhang, Q., Guo, Y., Sun, X., Xia, L., Darwish, I.A., 2024. Dual-  
932 aspect attention spatial-spectral transformer and hyperspectral imaging:  
933 A novel approach to detecting aspergillus flavus contamination in peanut  
934 kernels. *Postharvest Biology and Technology* 213, 112960. doi:10.1016/  
935 j.postharvbio.2024.112960.
- 936 Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., Ye, Q., Liu,  
937 T.Y., 2017. LightGBM: A highly efficient gradient boosting decision tree,  
938 in: *Advances in Neural Information Processing Systems*, pp. 3149–3157.  
939 doi:10.5555/3294996.3295074.
- 940 Khodabakhshian, R., Emadi, B., 2017. Application of vis/snir hyperspectral  
941 imaging in ripeness classification of pear. *International Journal of Food*  
942 *Properties* 20, S3149–S3163. doi:10.1080/10942912.2017.1354022.

- 943 Lanke, N., Chandak, M., 2025. Recent trends in deep learning and hyper-  
944 spectral imaging for fruit quality analysis: an overview. *Journal of Optics*  
945 doi:10.1007/s12596-025-02881-7.
- 946 Liu, D., Zhang, H., Lv, F., Tao, Y., Zhu, L., 2024. Combining transfer  
947 learning and hyperspectral imaging to identify bruises of pears across dif-  
948 ferent bruise types. *Journal of Food Science* 89, 2597–2610. doi:10.1111/  
949 1750-3841.17050.
- 950 Lu, B., Dao, P.D., Liu, J., He, Y., Shang, J., 2020. Recent advances of  
951 hyperspectral imaging technology and applications in agriculture. *Remote*  
952 *Sensing* 12, 2659. doi:10.3390/rs12162659.
- 953 Lu, Y., Huang, Y., Lu, R., 2017. Innovative hyperspectral imaging-based  
954 techniques for quality evaluation of fruits and vegetables: A review. *Ap-  
955 plied Sciences* 7, 189. doi:10.3390/app7020189.
- 956 Lundberg, S.M., Lee, S.I., 2017. A unified approach to interpreting model  
957 predictions, in: *Advances in Neural Information Processing Systems*, pp.  
958 4766–4777. doi:10.5555/3295222.3295230.
- 959 Min, D., Zhao, J., Bodner, G., Ali, M., Li, F., Zhang, X., Rewald, B.,  
960 2023. Early decay detection in fruit by hyperspectral imaging—principles  
961 and application potential. *Food Control* 152, 109830. doi:10.1016/j.  
962 foodcont.2023.109830.
- 963 Nagasubramanian, K., Jones, S., Sarkar, S., Singh, A.K., Singh, A., Gana-  
964 pathysubramanian, B., 2018. Hyperspectral band selection using genetic

965 algorithm and support vector machines for early identification of char-  
966 coal rot disease in soybean stems. *Plant Methods* 14, 86. doi:10.1186/  
967 s13007-018-0349-9.

968 Olisah, C.C., Trehella, B., Li, B., Smith, M.L., Winstone, B., Whitfield,  
969 E.C., Fernández, F.F., Duncalfe, H., 2024. Convolutional neural net-  
970 work ensemble learning for hyperspectral imaging-based blackberry fruit  
971 ripeness detection in uncontrolled farm environment. *Engineering Appli-  
972 cations of Artificial Intelligence* 132, 107945. doi:10.1016/j.engappai.  
973 2024.107945.

974 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel,  
975 O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., et al., 2011.  
976 Scikit-learn: Machine learning in Python. *Journal of Machine Learning  
977 Research* 12, 2825–2830. doi:10.5555/1953048.2078195.

978 Ribeiro, M.T., Singh, S., Guestrin, C., 2016. “why should I trust you?”: Ex-  
979 plaining the predictions of any classifier, in: *Proceedings of the 22nd ACM  
980 SIGKDD International Conference on Knowledge Discovery and Data Min-  
981 ing*, ACM. pp. 1135–1144. doi:10.1145/2939672.2939778.

982 Schmidt, V., Goyal, K., Joshi, A., Feld, B., Conell, L., et al., 2021. Code-  
983 Carbon: Estimate and track carbon emissions from machine learning com-  
984 puting. URL: <https://github.com/mlco2/codecarbon>, doi:10.5281/  
985 zenodo.4658424. software.

986 Su, Z., Zhang, C., Yan, T., Zhu, J., Zeng, Y., Lu, X., Gao, P., Feng, L.,  
987 He, L., Fan, L., 2021. Application of hyperspectral imaging for maturity

- 988 and soluble solids content determination of strawberry with deep learning  
989 approaches. *Frontiers in Plant Science* 12, 736334. doi:10.3389/fpls.  
990 2021.736334.
- 991 Varga, L.A., Makowski, J., Zell, A., 2021. Measuring the ripeness of fruit  
992 with hyperspectral imaging and deep learning, in: 2021 International Joint  
993 Conference on Neural Networks (IJCNN), IEEE. pp. 1–8. doi:10.1109/  
994 IJCNN52387.2021.9533728.
- 995 Vignati, S., Tugnolo, A., Giovenzana, V., Pampuri, A., Casson, A., Guidetti,  
996 R., Beghi, R., 2023. Hyperspectral imaging for fresh-cut fruit and vegetable  
997 quality assessment: Basic concepts and applications. *Applied Sciences* 13,  
998 9740. doi:10.3390/app13179740.
- 999 Wieme, J., Mollazade, K., Malounas, I., Zude-Sasse, M., Zhao, M., Gowen,  
1000 A., Argyropoulos, D., Fountas, S., Van Beek, J., 2022. Application of  
1001 hyperspectral imaging systems and artificial intelligence for quality assess-  
1002 ment of fruit, vegetables and mushrooms: A review. *Biosystems Engineer-*  
1003 *ing* 222, 156–176. doi:10.1016/j.biosystemseng.2022.07.013.
- 1004 Yang, C., Guo, Z., Fernandes Barbin, D., Dai, Z., Watson, N., Povey, M.,  
1005 Zou, X., 2025. Hyperspectral imaging and deep learning for quality and  
1006 safety inspection of fruits and vegetables: A review. *Journal of Agricultural*  
1007 *and Food Chemistry* 73, 10019–10035. doi:10.1021/acs.jafc.4c11492.