



**HAL**  
open science

# **cocoaDeep: a preliminary study of the performance sensitivity to datasets of Faster RCNN, YOLO and transformer networks for cocoa pod detection**

Philippe Borianne, Frédéric Théveny, Llorenç Cabrera-Bosquet, Sabine-Karen Lammoglia

## ► To cite this version:

Philippe Borianne, Frédéric Théveny, Llorenç Cabrera-Bosquet, Sabine-Karen Lammoglia. cocoaDeep: a preliminary study of the performance sensitivity to datasets of Faster RCNN, YOLO and transformer networks for cocoa pod detection. 2025. <hal-05339226>

**HAL Id: hal-05339226**

**<https://hal.sorbonne-universite.fr/hal-05339226v1>**

Preprint submitted on 30 Oct 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire HAL, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons CC BY-ND 4.0 - Attribution - No Derivative Works - International License

1 **cocoaDeep: a preliminary study of the performance sensitivity to datasets of**  
2 **Faster RCNN, YOLO and transformer networks for cocoa pod detection.**

3 Philippe Borianne<sup>1,2</sup>, Frédéric Théveny<sup>1,2</sup>, Llorenç Cabrera-Bosquet<sup>3</sup>, Sabine-Karen Lammoglia<sup>4,5</sup>.

4 <sup>1</sup>CIRAD, UMR AMAP, F-34398 Montpellier, France; <sup>2</sup>AMAP, Univ Montpellier, CIRAD, CNRS, INRAE, IRD,  
5 Montpellier, France; <sup>3</sup>LEPSE, Univ Montpellier, INRAE, Montpellier Institut Agro, Montpellier, France;  
6 <sup>4</sup>UFR Biosciences, Univ Félix Houphouët-Boigny, Abidjan, Côte d'Ivoire; <sup>5</sup>CIRAD, UMR ABSYS,  
7 Montpellier, France.

8 **Keywords**

9 cocoa pod detection, low-resolution images, neural networks, YOLO, Transformers, GPU  
10 times.

11 **Highlights**

- 12 • Pod detection on low-resolution cocoa tree images
- 13 • Comparative study of R-CNN, YOLO and Transformer networks
- 14 • True positives, false negatives and false positives
- 15 • Carbon footprint assessment
- 16 • Cocoa yield assessment

17

18

## 19 Abstract

20 Farmers must be able to estimate their crop yields at various growth stages for effective  
21 management of their farms and to enable them to interact with cooperatives or traders as early  
22 as possible. Here we developed an AI-based cocoa pod detection method using low resolution  
23 colour images of cocoa trees on farms in Côte d'Ivoire.

24 We compared Nano and eXtra-large architectures of six neural networks, including Faster  
25 RCNN, Baidu's Real-Time Detection Transformer (RTDetr), DETR-ResNet Vision  
26 Transformer, YOLOv5, YOLOv8 and YOLOv11.

27 These networks were trained with 7,850 annotated cocoa pods on 400 low resolution images,  
28 and validated in two independent datasets: a 42 low resolution images containing 990  
29 annotated pods, and a 100 low resolution images containing 2,400 annotated pods.

30 Unexpectedly, the Nano architectures of the YOLOv8 and YOLOv11 networks outperformed  
31 those of the RTDetr networks by ~2% and of the other networks by >5%, with an F1-score of  
32 77% on all images and up to 90% on foreground trees. The dominance of Nano architectures  
33 suggests that the eXtra-large architectures, which contain 20-30-times more neurons, may not  
34 have been fully trained due to insufficient data, thereby limiting the objectivity of inter-network  
35 comparisons. Finally, although the average detection performance of RTDetr for cocoa pods  
36 was only 2% lower than that of the YOLO8 network, it was definitively excluded from the  
37 candidate models because its per-image processing time was 15–20% higher than that of  
38 YOLOv8 and YOLOv11. However, with a performance sensitivity to data of less than 0.5%,  
39 YOLOv8 Nano became the best option.

## 40 Introduction

41 Estimating the yield of a crop at its various growth stages is essential for decision making on  
42 disease management, harvesting, storage, transport and marketing. Regarding cocoa  
43 (*Theobroma cacao*), a novel deep learning approach was recently presented for predicting  
44 cocoa yield using a recurrent neural network that combined spatiotemporal climatic data and  
45 statistical data on cocoa production in southwestern Nigeria [1]. However, due to a lack of  
46 suitable data, this model could not be transposed to Côte d'Ivoire, which alone accounts for  
47 >44% of world cocoa production [2], despite the fact that cocoa farming represents a major  
48 socioeconomic challenge for this country [3]. Yield estimation in that country is still based on  
49 manual counts, which are time-consuming and often hampered by major counting errors. This  
50 situation warrants the design of a low-cost, truly operational solution that could be integrated  
51 with current practices, while combining real-time identification of cocoa pods at different growth  
52 stages in natural environments with a yield estimation model based on pod counts—this would  
53 be essential for cocoa research and the cocoa industry. The main objective of this study was  
54 to propose a simple reliable AI-based method for image detection and counting of cocoa pods  
55 on low-resolution RGB mobile phone images.

56  
57 In recent years, several studies have described efficient computer vision systems for fruit  
58 detection and yield estimation [4]. Fruit detection involves finding instances of fruit in an image.  
59 Computer vision systems for fruit detection involve simple pixel segmentation in terms of  
60 density or colour or more advanced machine learning methods based on a combination of  
61 colour, shape and texture features computed on images sometimes acquired by multiple or  
62 multiband sensors. Advanced machine vision systems for yield estimation use deep learning  
63 algorithms for object detection: these are being used to an increasing extent in machine vision  
64 systems for yield estimation. Among deep learning algorithms, convolutional neural networks  
65 (CNNs) have highly efficient conventional detection performance [5]. A growing number of  
66 studies have used neural networks for fruit yield estimation [6], and/or fruit detection [7-8]

67 and/or fruit cultivar identification [9-10], while relying especially on Faster R-CNNs [11-13] or  
68 YOLO [14].

69 Numerous studies have compared diverse neural networks in terms of deep mechanisms of  
70 different families [15], as well as intra-family variations [16-17]. Already in 2020, a comparative  
71 study of different *in situ* fruit detection networks for apples, mangoes and oranges [18] placed  
72 Faster R-CNN in third position for all fruits, with an average accuracy ~2% lower than the best  
73 results obtained using an improved Faster R-CNN, and in second position for mango fruit  
74 detection, with a difference of 0.2% compared to the best results obtained with YOLOv3. More  
75 recently, a study in which a CNN model was compared to YOLO [19] for tomato detection  
76 revealed significantly superior YOLO results. A comparison between a detection transformer  
77 and YOLOv8 [20] for orange and sweet orange detection showed that the performance of  
78 these two models was relatively equivalent.

79 With regard to cocoa trees, most of the published studies we found concerned the detection  
80 of cocoa pod diseases [21-22] or ripeness [23]. These studies were based on images of cocoa  
81 pods, not of cocoa trees. The few studies involving plant pathology detection on cocoa images  
82 compared different convolutional neural networks and transformers, with scores ranging from  
83 80 to 90% [24]. [25] compared the two detection networks, i.e. U-Net and a fully convolutional  
84 network (FCN), where the latter performed much better, with a score of ~94%, i.e. nearly 2%  
85 higher than U-Net. Otherwise, [26] carried out a comparative study of different versions of  
86 Faster R-CNN and YOLOv5 detection networks, with the YOLOv5X, Faster R50FPN3x and  
87 R101C43x architectures obtaining almost identical scores of 95%.

88 Here we presented a comparative study of different neural network architectures for pod  
89 detection on colour images of cocoa trees. In particular, we compared Faster-RCNN [27],  
90 YOLOv5 [28], YOLOv8 [29], YOLOv11 [30], DETR-ResNet Vision Transformer [31] and the  
91 Real-Time Detection Transformer [32]. These networks were selected for three main reasons:  
92 (1) Faster RCNN is the original architecture that has historically been implemented for fruit  
93 detection, (2) the YOLO family has benefitted from continuous advances [33], and (3)

94 transformer-based models are emerging as successors to CNNs [34]. The simplest and most  
95 complex architectures were deployed and the F1-score of each network was compared to  
96 assess the extent of precision loss, particularly in cases where the cocoa pod detector could  
97 be directly embedded in the smartphone or tablet acquiring the images. CNN ResNet50 and  
98 ResNet101 architectures [35] were compared with YOLO's specific architectures. In addition,  
99 the complexity of objectively estimating the digital processing carbon footprint [36] was  
100 addressed to generate all elements needed to assess the efficiency of future operational  
101 solution. This exploratory study was part of a broader context focusing on early yield  
102 estimation. The yield estimate for a cocoa plot is often obtained by extrapolating the production  
103 estimate for a few cocoa trees chosen at random from the plot; the production of each of these  
104 cocoa trees is itself estimated based solely on the cocoa pods visible and counted in a  
105 photograph of the tree. The networks should therefore regularly process datasets from different  
106 cocoa tree plots. This study did not seek to investigate the optimal performance of each neural  
107 network trained and tested on a specific data set; it focused specifically on each trained  
108 network (performance) sensitivity to datasets. The study deliberately did not focus on the  
109 results of tests carried out during the network training sessions; rather, it focused on the results  
110 obtained on different validation datasets that were similar and representative of the data given  
111 to networks during early plot yield estimates. The study also aimed to assess which models  
112 for detecting cocoa pods could be installed on smartphones in order to eventually offer field  
113 solutions.

114 The first part of this article focuses on the data and their origins, the characteristics of the  
115 studied networks and the principles applied for estimating the carbon footprint. The second  
116 part presents the values of the indicators used and the main results of the comparative study.  
117 Finally, the third part contains a general discussion on the results, their relevance and future  
118 prospects.

## 119 Materials and methods

### 120 Material

121 Cocoa (*Theobroma cacao*) is a small evergreen tree of the Sterculiaceae family. It produces  
122 edible cocoa beans with different flavours depending on the cocoa tree variety, and cocoa is  
123 the main ingredient in chocolate making [37]. These beans develop in the fruit (pods) of the  
124 cocoa tree. Cocoa pods are elongated and resemble a fairly rounded cucumber. They measure  
125 15-25 cm in length and 6-15 cm in diameter, while weighing 300-500 g depending on the  
126 variety. These pods are usually reddish-yellow in colour when ripe (see Fig. 1). Because of the  
127 large mass of seeds contained in cocoa pods, the pod surface is covered with numerous small  
128 bumps, but also marked by around 10 relatively deep longitudinal grooves.



129

130 *Figure 1: Forastero cocoa orchards, Abidjan district in western-central Côte d'Ivoire.*

131 Cocoa growing is not mechanised in Côte d'Ivoire [38] and cocoa farms are fairly dense  
132 'natural' orchards (see Fig. 1). Forastero is the most widely grown variety, which accounts for  
133 79% of cocoa production in the country, whereas the Trinitario variety accounts for 20%. Criollo  
134 is the rarest cocoa variety, and is currently the least cultivated, accounting for just 1% of all  
135 harvested cocoa [39]. Our study focused exclusively on Forastero cocoa.

136 Image acquisition

137 The cocoa image dataset was acquired manually by an operator in cocoa crop fields in Côte  
138 d'Ivoire (July 2019). The images were captured using the main camera of a mobile Samsung  
139 Galaxy S10 SM-G973F phone in uncontrolled field conditions. No specific protocols were  
140 followed for image acquisition, and parameters such as shooting distances and focal lengths  
141 were not specified in advance. Consequently, images in the dataset exhibited variations in  
142 depth of field and shooting angles. The lighting conditions and fruit occlusion were not  
143 controlled, and the images were taken without flash. Images of different random trees were  
144 captured, and the distance between the trees and the camera sensor was not consistently  
145 recorded. This lack of a predefined image acquisition protocol and the use of a handheld device  
146 in real-field conditions accounted for variability in the dataset for some factors such as  
147 illumination and object scale. The resulting dataset consisted of 500 colour images with 4,600  
148 x 3,456 pixel resolution. These images were then deliberately resized to a low resolution of  
149 1,008 x 756 pixels in order to reduce processing times on both computer servers and  
150 smartphones.

151

152 Image annotation

153 Expert image annotation is very crucial as it enables the creation of datasets to train neural  
154 networks, and of validation datasets to assess their performance. This is a tedious and time-  
155 consuming task that often requires substantial resources, as the annotation quantity needs to  
156 be as high as possible to be able to conduct highly efficient performance studies. The authors  
157 deployed the cocoa-fruit-counting project on Zooniverse [40], the world's largest participatory  
158 research platform [41]. Initially developed for the study of constellation images, the platform  
159 allows millions of volunteers around the world to be called upon to enhance research project  
160 resources, which was particularly helpful here for the image annotation stages.

161



Figure 2 : Zooniverse annotation. The area of interest outlined by the red polygonal line indicated the tree in the foreground. All the cocoa pods visible in the image were outlined by ellipses shown in blue: their adjustment was based on translations, rotations and resizing of a canonical ellipse.

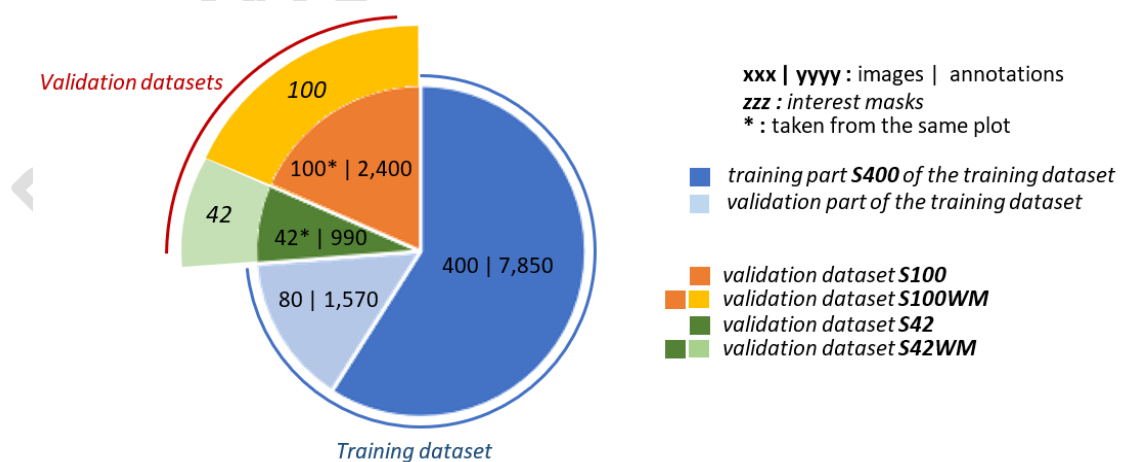
162

163 The Zooniverse cocoa pod annotation process consisted of circumscribing each cocoa pod  
 164 with an ellipse (in blue in Figure 2). As this research was part of a yield estimation study, each  
 165 Zooniverse operator (trained expert) was asked to outline the foreground tree on each image  
 166 with a polygon (in red in Fig. 2). Once completed, the annotations were validated and corrected  
 167 where necessary via the ImageJ platform [42].

168 Training and validation datasets

169 The annotated data was divided into disjoint training and validation datasets that did not share  
 170 any data in common. The labelled annotations were deduced from the Zooniverse ellipses and  
 171 consisted of rectangular bounding boxes with dimensions ranging from 10 to 100 pixels per  
 172 side. These bounding boxes were drawn to as accurately as possible enclose the visible  
 173 portions of the pods on the images. The training dataset included ~7,850 annotations

174 representing diverse visual cocoa pod features in terms of shape, colour, sunlight conditions  
 175 or occlusions on >400 images (S400) from different plots: additional dataset has been added  
 176 to the training dataset to estimate the network's performance during the various training loops.  
 177 Two distinct un breakable validation datasets were created with the remaining annotated data  
 178 to assess variability in the trained networks with the aim of ensuring that the training and  
 179 validation datasets would be representative of the overall data distribution. The first validation  
 180 dataset (S42) consisted of 42 images that included 990 annotated pods. The second validation  
 181 dataset (S100) contained 100 images with a total of 2,400 annotated pods. Each validation  
 182 dataset consisted of images from a single plot, following the recommended data acquisition  
 183 protocol for early yield estimation. The difference in dataset size is solely due to the difference  
 184 in plot size. The plots were selected from the same geographical region, with the same  
 185 cultivation practices and the same varieties of cocoa trees. These two sets were therefore  
 186 considered similar, the idea being to assess the sensitivity of the detection networks to the  
 187 data. Furthermore, the polygons outlining the foreground tree were used to generate binary  
 188 interest masks. These masks were black-and-white images on which the region defined by the  
 189 polygonal boundary of the area of interest cocoa tree was rendered in white, while the  
 190 remaining area was black. Added to datasets S42 and S100 respectively, they defined  
 191 datasets S42WM and S100WM (see Fig. 3).



197 *Figure 3: training and validation datasets. The annotated images are divided into three non-overlapping*  
 198 *datasets: the S400 training dataset, consisting of 400 images from different plots, and the S42 and S100*  
*validation datasets, consisting of 42 and 100 images from the same plots, respectively.*

199 Neuronal networks

200 The study was focused on Faster-R-CNN, YOLOv5, YOLOv8, YOLOv11, DETR-ResNet  
201 Vision Transformer and Baidu's Real-Time-Detector networks. While the Region-based  
202 Convolutional Network (R-CNN), You Only Look Once (YOLO) and Transformers are object  
203 detection networks, their respective mechanisms differ. R-CNNs start by finding interesting  
204 parts of the image, then examine these parts more closely to determine what they contain.  
205 YOLO networks grid the space via successive convolutions and then locate and identify objects  
206 of interest. This model first emerged in 2015, but it was not until YOLOv5 (2021) was designed  
207 that the model became easily accessible and sufficiently optimised for acceptable detection.  
208 YOLOv8 (2023) was then the first model that offered truly significant features and  
209 improvements for enhanced detection performance, flexibility and efficiency. YOLOv11 (2024)  
210 is the latest iteration in the “classic” series of real-time object detectors. This version includes  
211 major improvements in the architecture and learning methods, so it is more versatile than its  
212 predecessors. Self-attentive Transformers models (2017) are conceptually based on the  
213 attention mechanism and take the context of the observations to be predicted into account.  
214 This type of model is particularly effective in translation tasks, where a sequence of words in  
215 one language is transformed into a sequence of different words in another language. Vision  
216 Transformers emerged in late 2020 and adapted the architecture of Transformers to process  
217 image content. These models are rapidly overtaking and replacing convolutional models for  
218 computer vision. The Real-Time Detection Transformer features significant optimisations for  
219 dealing with multi-scale features—this involved decoupling intra-scale interaction and inter-  
220 scale fusion, and making the model highly adaptable, thereby enabling flexible inference speed  
221 adjustment.

222 Yet beyond their differences, these networks share common features. They are all based on a  
223 neural architecture i.e., an arrangement of successive layers containing unitary computing  
224 elements called neurons. The neurons are all weighted and biased during training in order to  
225 be suitable for the detection of objects of interest. The total number of parameters is a linear

226 combination of the weights, biases, inputs and outputs of each network layer. The Faster-R-  
227 CNN Nano ResNet50 architecture (FRCNNN) has around 45 million parameters, the Faster-  
228 R-CNN eXtra-large ResNet101 architecture (FRCNNX) around 55 million, the YOLOv5 Nano  
229 architecture (YOLOv5N) around 1.8 million, the YOLOv5 eXtra-large architecture (YOLOv5X)  
230 around 86 million, the YOLOv8 Nano architecture (YOLOv8N) around 3 million, the YOLOv8  
231 eXtra-large architecture (YOLOv8X) around 86 million, the YOLOv11 Nano architecture  
232 (YOLOv11N) around 2.5 million, the YOLOv11 eXtra-large architecture (YOLOv11X) around  
233 57 million, the DETR-ResNet Vision-Transformer Nano ResNet50 architecture (ResNetDetrN)  
234 around 26 million, the DETR-ResNet Vision-Transformer eXtra-large ResNet101 architecture  
235 (ResNetDetrX) around 86 million, the Baidu's Real-Time-Detector Large architecture  
236 (RTDetrL) around 32 million and the Baidu's Real-Time-Detector eXtra-large architecture  
237 (RTDetrX) around 66 million. Furthermore, the networks can only structurally process images  
238 of fairly low resolution, i.e. generally between 300 and 1,096 pixels. When images are larger  
239 than the resolution limits accepted by the network, the latter automatically resizes the native  
240 image, although this sometimes alters the image content.

#### 241 Image tiling

242 Small objects of interest may disappear or become indistinguishable after resizing the native  
243 images. To avoid this major inconvenience, it is advisable to divide the native image into  
244 thumbnails, or so-called tiles, of a resolution that can be processed by the networks, with small  
245 objects of interest remaining visible and detectable. A tiling strategy with tile overlap was  
246 adopted to ensure that each object of interest would be fully contained within at least one tile.  
247 The detection networks therefore provided predictions, i.e. bounding boxes locating cocoa  
248 pods in each tile. Post-processing was applied to transfer the cocoa pod detections obtained  
249 in each tile to the original image. This involved repositioning the predicted bounding boxes  
250 from the individual tiles on the coordinate system of the full image. Due to tile overlapping,  
251 duplicate detections of the same cocoa pod could occur across adjacent tiles, so a process  
252 was implemented to delete duplicate predictions.

253 Carbon footprint

254 A computing server typically comprises two distinct but complementary resources: the service  
255 resource that manages the web application interfacing the user and the computations, and the  
256 compute resource alternately mobilising the computer's central processing unit (CPU) cores  
257 and the graphical processing unit (GPU) cores of the computer graphics card. The GPU  
258 manages the storage of temporary data generated during computations. The calculation server  
259 is also often associated with a server for storing input data and results produced for periods of  
260 varying length, as generally specified in the calculation server guidelines. The whole system is  
261 housed in a data centre equipped with power supply and cooling systems.

262 Estimating the carbon footprint may be very complicated. Various factors must be taken into  
263 account, including the energy consumption of the hardware components (CPU, GPU, memory,  
264 storage) during the different operational phases (training, validation, inference), the efficiency  
265 of the power supply and cooling systems in the data centre, and the source of the electricity  
266 powering the infrastructure. While detailed carbon footprint analysis is beyond the scope of  
267 this study, information on the CPU and GPU processing times and the temporary data volumes  
268 were provided as carbon footprint indicators.

269 Design of experiments (DOE)

270 Each network was trained on the S400 training dataset and its performance was evaluated on  
271 the two disjoint validation datasets S42 and S100, with or without interest masks. The  
272 validation of datasets with interest masks consisted of pairing only expert annotations and  
273 network predictions whose geometric centre of the bounding box belonged to the area of  
274 interest. Benchmark performances were provided based on the S400 training set. The  
275 hyperparameters for each network were left at the default values proposed by their respective  
276 authors, who defined the optimal settings for their model. The batch size, which defines the  
277 number of tiles processed simultaneously by the network, was arbitrarily set at 32 based on  
278 the memory of the graphics card used. The metric chosen during training was the F1-score for  
279 convenience: it is a well-known indicator used by the agronomy community. The same

280 minimalist data enrichment was applied to each network: it consisted solely of systematic  
281 vertical inversion of the images on the one hand, and systematic 90° rotations on the other, to  
282 'show' the network of horizontal dents what happens when images accidentally switch to  
283 landscape mode due to the orientation of smartphones but are processed in portrait mode.  
284 The maximum number of epochs i.e., the maximum number of times the network would see  
285 all the data, was set at 1,200 with a patience of 100 epochs, indicating the number of epochs  
286 to wait without improvement in validation metrics before stopping training early. Thus, the  
287 effective training of YOLO took place between 850 and 930 epochs depending on the model,  
288 unlike RCNN. Performance evaluation was conducted at two levels, i.e. at the image level,  
289 where all predictions within the entire image were considered, and at the foreground tree level,  
290 where the evaluations were focused on or limited to the masks of interest defined per image.  
291 This allowed for a detailed analysis of the detection accuracy, both generally and within specific  
292 regions.

293 The normalized rate of geometric overlap between two bounding boxes was used to measure  
294 the total or partial superposition of objects. The intersection on union (IoU) of objects was  
295 defined by the ratio of the intersection to the union of the two bounding boxes [43]. This latter  
296 indicator was applied for the deletion of cocoa pods predicted several times by the networks  
297 due to tile overlap bands and for matching expert annotations with network predictions. The  
298 value of this ratio was compared with a pre-determined threshold and used to identify  
299 overlapping boxes considered to represent the same cocoa pod. The object detection  
300 performance was estimated according to the F1-score [44], which is a standard comparative  
301 test between the image truth (expert annotations) and the network predictions. The F1-score  
302 was the harmonic mean between the precision and recall, i.e. two statistical indicators giving  
303 the contribution of 'false positives' and 'false negatives' respectively to the network's overall  
304 detection error: true-positives were cocoa pods annotated by the expert and correctly detected  
305 by the network, false-negatives were cocoa pods annotated by the expert but not detected by

306 the network, and false-positives were also cocoa pods detected by the network but not  
307 considered as such by the expert.

308 The performance indicators were assessed by cocoa pod size class. Based on the dimensions  
309 of the ground truth bounding boxes, cocoa pods framed by boxes smaller than 16x16 pixels  
310 were considered 'small', cocoa pods framed by boxes between 16x16 and 32x32 pixels were  
311 considered 'medium', and cocoa pods that did not fit into either of the above categories were  
312 considered 'large'.

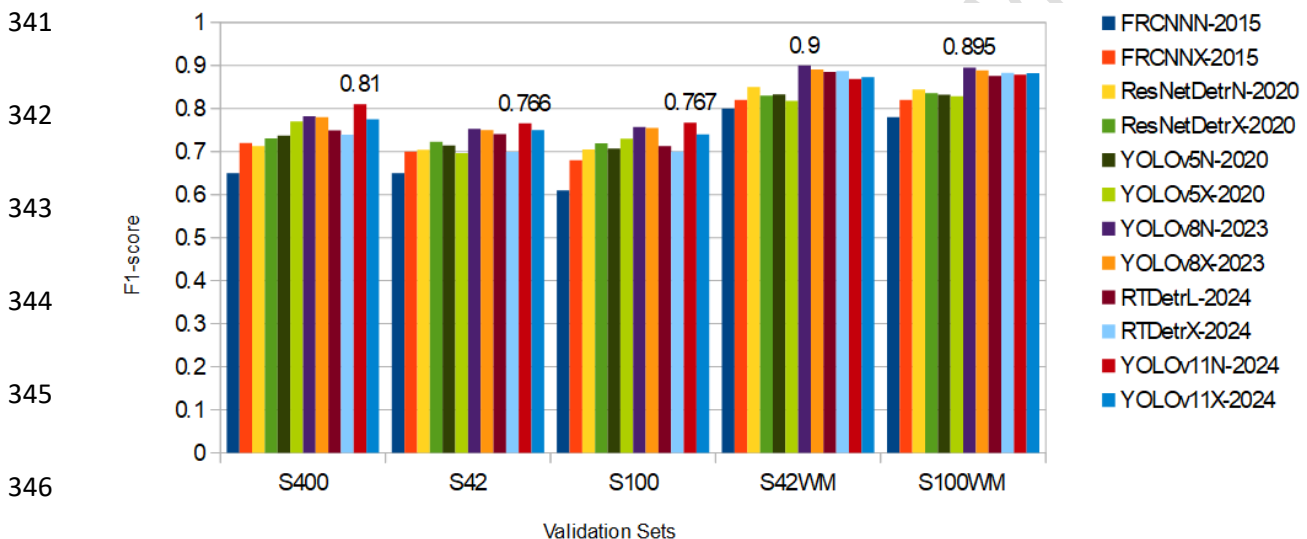
313 All the experiments were carried out on the AgroDeep platform [45], thereby ensuring that the  
314 technical environments were relatively similar between the different tests. Python scripts for  
315 pre-processing the data, post-processing the results and inferring the neural networks were  
316 encapsulated in Singularity containers [46], which provided an optimal operating environment  
317 (operating system, additional calculation packages, CUDA drivers, etc.) for each network in  
318 the study. All training and validation procedures were carried out on the same computing  
319 server, which was exclusively dedicated to the study, and equipped with a 24GB Nvidia Quadro  
320 RTX6000 graphics card, 4,600 CUDA cores and 570 Tensor cores.

321 The tile size was set at 640x640 for training all the networks and at 640x640 with a 200 pixel  
322 overlap for validating all the networks. The threshold for removing multiple predictions was set  
323 at 0.25. A minimum IoU threshold of 0.7 was set to determine a successful match between the  
324 network predictions and the expert ground truth annotations.

## 325 Results

326 The S400, S42 and S100 datasets illustrated the network performances at the entire image  
327 scale, while the S42WM and S100WM datasets illustrated these performances solely at the  
328 foreground tree scale. Unsurprisingly, the most recent networks (published as of 2023) gave  
329 better results than the older networks (see Fig. 4). The top three were clearly RTDetr, YOLOv8  
330 and YOLOv11. With an F1-score of ~77% for S42 and S100, the Nano architecture of  
331 YOLOv11 gave markedly better results than YOLOv8 and RTDetr, with F1-scores ranging from

332 75.3% to 75.7%. The trend was, however, reversed for the S42WM and S100WM datasets,  
 333 i.e. with an F1-score of ~90%, the YOLOv8 Nano architecture was around 2% better than  
 334 YOLOv11 and RTDetr. With a maximum difference of 3%, the F1-scores of the Nano and  
 335 eXtra-large architectures of the YOLOv8 and YOLOv11 networks were very close, or even  
 336 equivalent, while being substantially higher than those of the RCNN, YOLOv5 and  
 337 transformers, with a difference of >5%, all architectures combined. Pre-experimentation on the  
 338 validation of the different networks trained on a set of 255 images containing 4,527 annotated  
 339 cocoa pods led to a systematic performance decrease of 5-15% on the S42, S100, S42WM  
 340 and S100WM datasets.



347 *Figure 4: summary of the performance of different networks. Performance was measured by the F1 score of each network*  
 348 *trained under similar conditions on each of the validation sets: the S400 set provided the reference performance of the*  
 349 *networks on their own training set.*

349 The histograms for S42 and S100 (resp. S42 WM and S100WM) were similar in shape and  
 350 value, thus indicating the relative stability of the response of these networks when assessed  
 351 with the different validation datasets.

352 The F1-scores obtained for S400 were 2-4% higher than those obtained for S42 and S100, as  
 353 expected since the S400 dataset was exclusively composed of the network training images.  
 354 This slight drop in performance observed for validation datasets containing data that the

355 networks had not encountered during their training highlighted the model's generalisation  
 356 potential.

357 Assessing the sensitivity of network performance to data sets was difficult with so little data  
 358 available. However, trends were already emerging when comparing network performance  
 359 obtained on the S42 and S100 validation datasets, first without and then with the use of interest  
 360 masks (see Fig. 4). The largest variation observed was 4% for Faster RCNN; this variation fell  
 361 below 1% for Detr-ResNet transformer, YOLOv5 and the large-RTDetr transformers, and fell  
 362 to less than 0.5% for the others. The candidate networks were therefore those that performed  
 363 close to 90% with a variation of less than or equal to 0.5%.

364 We therefore focused our study on the top three networks with high performance on validation  
 365 datasets and low performance variations between validation datasets.

366 The F1-scores of the Nano and eXtra-large architectures of the YOLOv8, YOLOv11 and  
 367 RTDetr networks by cocoa pod size class (see Tab.1) were used to refine the above overall  
 368 results, with the lowest scores indicated in red, the medium scores in yellow and the highest  
 369 scores in green. The lowest scores were obtained for small probably embryonic cocoa pods in  
 370 the foreground or mature fruit in the background, while the highest scores were obtained for  
 371 large probably mature cocoa pods on foreground trees. This initial result confirmed that the  
 372 networks had more difficulty in detecting small background fruit than large foreground fruit.

373 *Table 1 : F1-scores in percentage of the Nano and eXtra-large architectures of the YOLOv8, YOLOv11 and RTDetr networks by*  
 374 *cocoa pod size class. Colour coding of the cells ranging from red for the lowest scores to green for the highest scores highlighted*  
 375 *the networks' poor ability to detect small cocoa pods in the background and their strong ability to detect large pods in the*  
 376 *foreground.*

size classes →	YOLOv8N			YOLOv8X			YOLOv11N			YOLOv11X			RTDetrL			RTDetrX		
	small	med.	large	small	med.	large	small	med.	large	small	med.	large	small	med.	large	small	med.	large
S400	46.97	77.08	92.54	50.04	77.78	92.73	50.98	82.16	92.97	43.95	78.68	91.58	45.23	72.32	91.3	46.13	73.51	93.46
S42	31.75	85.62	94.92	31.6	85.8	95.1	33.91	80.98	93.26	35.73	84.6	91	37.63	81.05	91.47	35.48	81.99	93.33
S100	45.51	83.5	95.05	49.06	83.25	94.5	48.58	83.65	94.9	39.3	83.76	93.31	43	77.71	92.05	43.65	80.46	94.22
S42WM	54.55	89.05	95.31	55.74	88.62	93.59	52.17	85.2	93.63	46.67	87.94	91.78	25	86.19	94.37	37.21	86.96	95.2
S100WM	45.71	88.44	95.69	47.9	88.12	95.36	42.38	86.34	95.44	41.33	88.46	94.09	30.93	85.3	94.1	41.1	87.2	94.1
av. NWM	41.41	82.07	94.17	43.57	82.28	94.11	44.49	82.26	93.71	39.66	82.35	91.96	41.95	77.03	91.61	41.75	78.65	93.67
av. WM	50.13	88.74	95.5	51.82	88.37	94.48	47.28	85.77	94.53	44	88.2	92.94	27.96	85.74	94.23	39.15	87.08	94.65

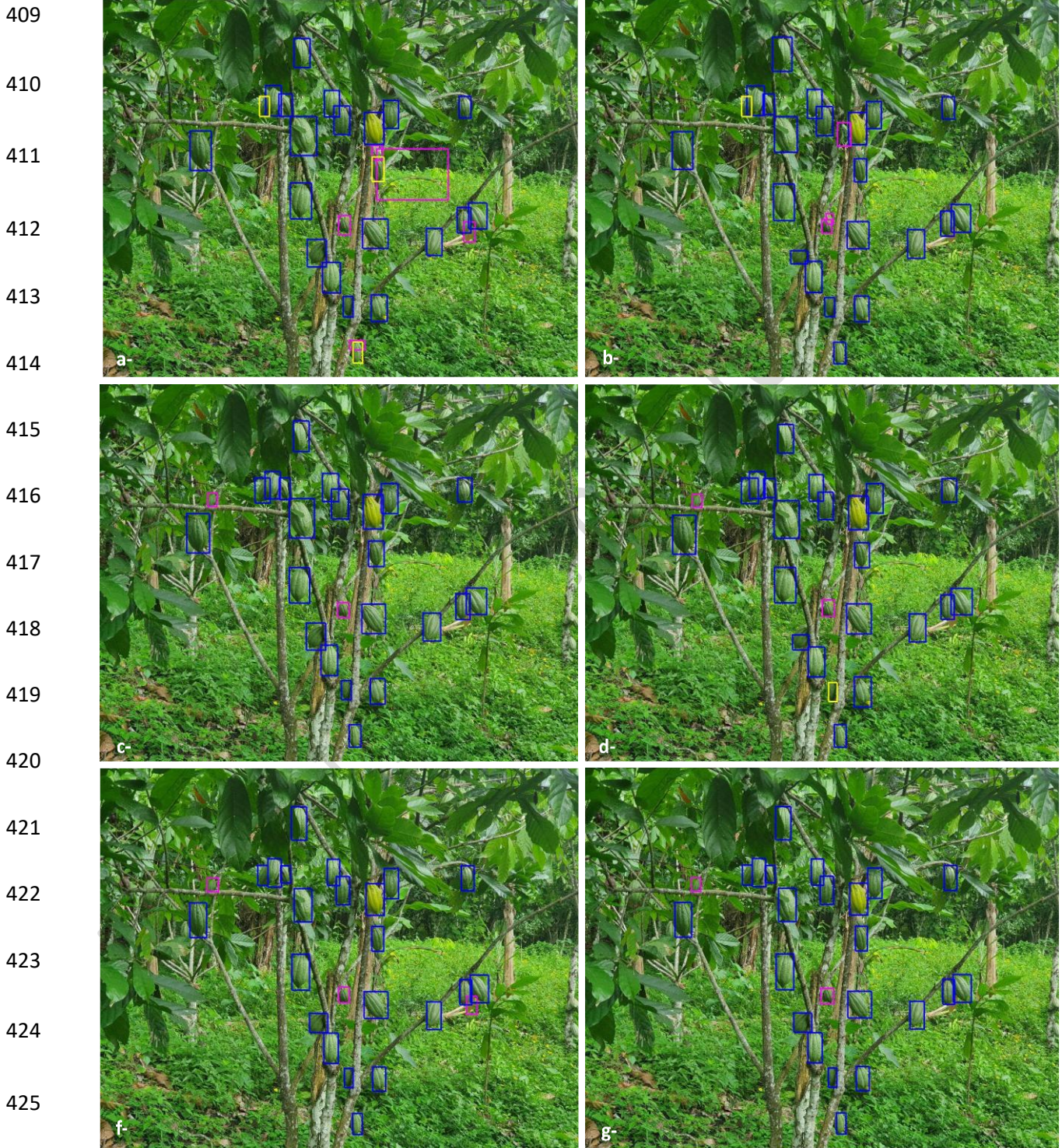
379 The penultimate row of the table gives the average scores per fruit size class for cocoa pod  
380 detection on whole images, while the last row gives the cocoa pod detection scores only in  
381 terms of the trees of interest. A slight improvement in the average F1-scores per network was  
382 observed when the cocoa pod detections were limited to foreground areas on the images of  
383 the trees of interest.

384 The networks could be ranked according to their respective performances solely based the  
385 analysis of the results obtained for the trees of interest. An average F1-score of 43% for small  
386 cocoa pod detection, but with a dispersion of ~9%, showed that the RTDetr architectures did  
387 not perform well in this size range. An F1-score of 87% for medium-sized pod detection with a  
388 dispersion of <1.5% illustrated the similarity of the responses of the three networks, although  
389 YOLOv8 had a slight advantage over YOLOv11 and RTDetr. With an F1-score of 94% for large  
390 cocoa pod detection with a dispersion of <1%, YOLOv8, YOLOv11 and RTDetr confirmed their  
391 ability to identify and count pods on foreground trees. With average F1-scores of 78% across  
392 all size classes, the eXtra-large and Nano architectures of the YOLOv8 network were ~3%  
393 better than YOLOv11 and ~6% better than RTDetr. But with an F1-score of 95.5% for the large  
394 fruit class alone, the YOLOv8 Nano architecture performed >1% better than all the other  
395 architectures.

396 A qualitative assessment was achieved by images showing the match between expert  
397 annotations and network predictions: the true-positives were represented by blue bounding  
398 boxes, the false-negatives by yellow bounding boxes and the false-positives by pink bounding  
399 boxes in the zones of interest on the foreground tree. Although not easy to appraise, the blue  
400 and pink boxes showed the location of the network predictions, which meant that there could  
401 potentially have been slight differences between networks.

402 On the foreground tree image in Figure 5, Nano YOLOv8 (c) and eXtra-large RTDetr (g)  
403 obtained the fewest errors, with 22 true-positives and two false-positives, for a local F1-score  
404 of 95.6%. All annotated cocoa pods were detected, and also the network detected two false-  
405 positives: the pink box in the centre of the image was a FALSE false-positive, i.e. a real pod

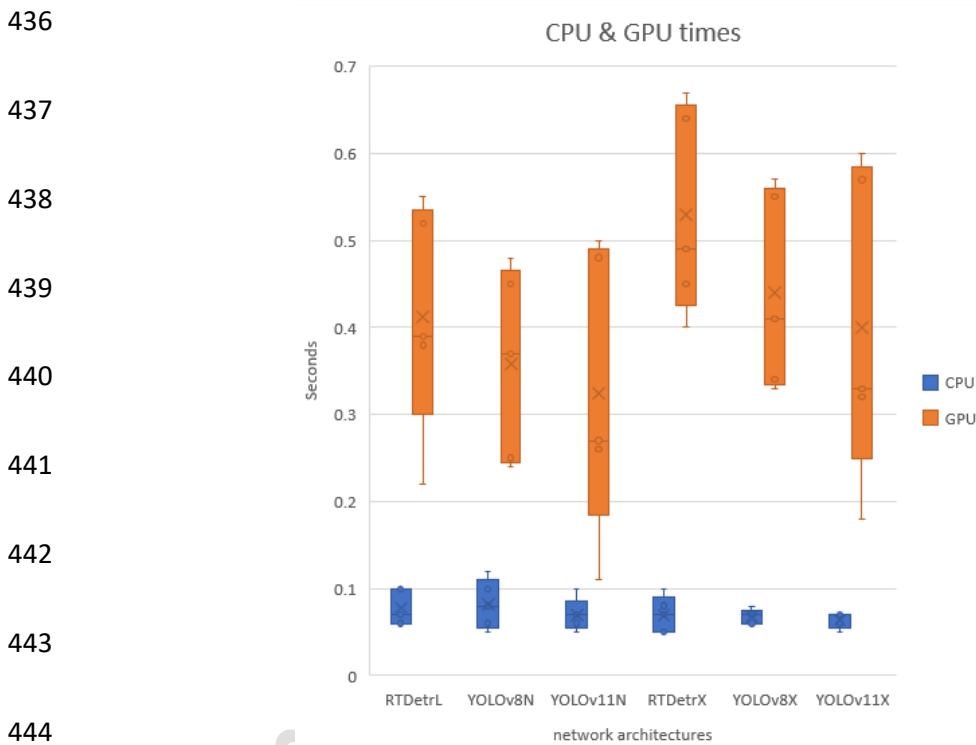
406 not annotated by the expert; the pink box furthest to the left was apparently related to a  
407 background cocoa pod captured by the network due to an insufficiently tight contour around  
408 the foreground tree.



426 *Figure 5: Qualitative visual assessment of prediction accuracy for YOLOv11 (top), YOLOv8 (middle) and RTDetr (bottom) in*  
427 *Nano (left) and eXtra-large (right) architectures. The prediction boxes are coloured according to the annotation/prediction*  
*appearing classes:*

428 Paradoxically, eXtra-large YOLOv8 (d) also missed a cocoa pod, despite having a more  
 429 complex architecture than Nano YOLOv8. The eXtra-large YOLOv11 architecture (b) seemed  
 430 to be more efficient than the Nano YOLOv11 architecture (a), which generated numerous  
 431 detection errors, although not as efficient as the YOLOv8 architecture. With three false-  
 432 positives and a local F1-score of 93.3%, Nano RTDetr (f) ranked numerically *ex aequo* with  
 433 networks presenting one false-negative and two false-positives (b,d), even the absence of  
 434 false-negatives was a good thing.

435 CPU and GPU times



445 *Figure 6: Mean CPU and GPU times for RTDetr, YOLO8 and YOLO11 validations per image. Average times were assessed solely on the S400, S100 and S42 validation datasets.*

446 The CPU and GPU runtimes for training and network validation were defined by differences  
 447 between the system times taken at the start and end of each key stage. The server used was  
 448 exclusively dedicated to the calculations so as to ensure inter-network comparability of  
 449 calculation times. Input data tiling, removal of multiple predictions, pairing of expert annotations  
 450 and network predictions and processing performance estimations were carried out via the  
 451 computer's CPU, whereas the learning and prediction of networks were carried out via the

452 computer's GPU. The CPU time to load data from the hard disk to the graphics card memory  
453 during GPU processing was overlooked.

454 The average unit times per type of computing resource, evaluated on the validation of the five  
455 batches S400, S100, S42, S100WM and S42WM, accounted for the resources consumed per  
456 network. The boxplots (see Fig. 6) summarise the datasets at six key figures, i.e. the minimum  
457 (lower horizontal line), the 25th percentile (lower circle), the median (horizontal line), the mean  
458 (cross), the 75th percentile (upper circle) and the maximum values (upper horizontal line),  
459 thereby facilitating the detection of outliers in the distribution of resource consumption values,  
460 which could indicate instances of particularly high or low resource usage. With a tiny dispersion  
461 of 0.0066 s, the average CPU times/image between the different networks were considered to  
462 be similar—a very logical result given that the input data pre-processing and output post-  
463 processing algorithms were rigorously identical. With a significant dispersion of 0.071 s, the  
464 GPU times illustrated the impact of the different complexities and mechanisms of the  
465 architectures tested. RTDetr transformer was clearly the slowest network, with an average  
466 GPU time/image of 0.41 s for its simplest architecture and 0.53 for its most complex  
467 architecture. YOLOv8 was the fastest network, with an average time of 0.32 s for its simplest  
468 architecture and 0.40 s for its most complex architecture. However, with the exception of the  
469 Nano architecture of the YOLOv8 network, the boxplots showed outlier distributions, with the  
470 median value tending towards the 25th percentile. For YOLOv8, the median tended towards  
471 the mean.

## 472 Discussion and prospects

473 With overall F1-scores <80% on all images, and barely 90% on foreground trees, cocoa pod  
474 detection was less successful than the detection of apples and sweet peppers [11], mangoes  
475 [12], oranges [14] and lemons [20], which had F1-scores of 92-98%. The score estimation  
476 depended on the validation dataset contents, which could likely be explained by the nature  
477 and repeatability of the training data. The number of annotations (or annotated images)

478 seemed very sufficient, yet the use of Zooniverse induced major bias. Since the same images  
479 were successively annotated by several users, the datasets included multiple annotations of  
480 the same cocoa pods, although there were no real differences in context. The 11,240  
481 annotations involved in the study accounted for <2,500 truly different cocoa pod contexts.  
482 There were not enough examples to constitute sufficiently representative training and  
483 validation datasets. The relationship between the training dataset size and the network  
484 performance has long been known [47], as was confirmed here by network performance  
485 differences of 5-15% when comparing training with a set of 225 annotated images and with a  
486 set of 400 annotated images.

487 Assessing the foreground tree on low-resolution images sometimes proved difficult, as the  
488 viewpoint could not be changed. Drawing a single polygon to outline the foreground tree proved  
489 restrictive, i.e. cocoa pods of background trees were visible in the polygon of interest  
490 circumscribing the foreground tree due to the low cocoa tree foliage density; the same problem  
491 also arose for cocoa pods on peripheral foreground tree branches located within the polygon  
492 outlining the tree of interest. However, when detected by the network, these specific pods were  
493 considered false positives and contributed to increasing the network's error rate. This  
494 observation highlighted a major difference between the approaches taken by the expert and  
495 the neural network: the expert mentally identified the branches of the tree of interest on which  
496 he had to annotate the cocoa pods, unlike the neural networks, which detected all the pods in  
497 the area of interest, regardless of the tree on which they were growing. The accuracy of the  
498 delineation of the area of interest was therefore crucial.

499 The defined acquisition protocol allowed shooting distances in the 3-12 m range to be fully in  
500 line with field operators' practices and constraints. A tree photographed at 3 m distance thus  
501 showed relatively large cocoa pods, whereas the same pods appeared small when  
502 photographed at 12 m distance. This flexible shooting range should enhance the detection  
503 network robustness and enable operators to implement a size class approach. Our cocoa pod  
504 size class assessment method partitioned the data and substantially improved the results, i.e.

505 boosting the F1 scores to ~94%, yet it could be even further improved by switching from pixel  
506 to metric dimensions. When using centimetric dimensions, the images were calibrated from  
507 reference patterns on the foreground tree or from image metadata as described by [47]. The  
508 cocoa pod detection networks were ranked according to two complementary criteria: the value  
509 of their respective performance and the low sensitivity to the data of this performance. This  
510 performance was assessed using the F1 score, widely used in thematic publications [48]; it  
511 had the dual advantage of being widely used and understood by agronomists in their own  
512 validation studies while incorporating detection excesses (false positives) and defects (false  
513 negatives). Direct interpretation of the inter-network classification results was therefore  
514 complex. Score differences >5% were considered relevant enough to validate the inter-network  
515 ranking, but differences <5% were not considered sufficiently relevant: we considered that  
516 these differences could be explained more by the datasets used than by the actual network  
517 capabilities. In this specific case, the performance of the networks could be considered  
518 equivalent, and the final choice of model to be implemented in an operational solution could  
519 be based on an additional criterion such as processing time. Beyond this ranking process, our  
520 study aimed to assess whether the drop in performance of the Nano architectures which, unlike  
521 the eXtra-large architectures, favoured speed over precision [30], was still acceptable for  
522 performance estimates. Against all expectations, the YOLOv8 and YOLOv11 Nano  
523 architectures achieved F1-scores that were a few hundredths higher than those of the eXtra-  
524 large YOLO and RTDetr architectures. Figure 5 clearly illustrated the difficulty of the eXtra-  
525 large YOLOv8 and YOLOv11 architectures in detecting all pods on the foreground tree, thereby  
526 suggesting that these networks had not been properly trained. The eXtra-large architectures  
527 had 20- to 30-times more parameters to weight than the Nano architectures. Indeed, they  
528 would have needed much more data to achieve training of equivalent quality [49], which would  
529 be necessary to be able to truly objectively compare the networks.

530 The performance sensitivity study was difficult to conclude due to the obvious lack of data. A  
531 dozen additional validation datasets from different cocoa plots would have been necessary;

532 however, two validation datasets were sufficient to identify the main trends and a few rules for  
533 selecting candidate neural networks. It is not unreasonable to exclude from the list those  
534 networks whose performance variations exceeded 0.5%.

535 The CPU times per image were stable, even though inter-network variations of  $<0.007$  s were  
536 observed. These variations did not markedly influence the time required to count cocoa pods  
537 in a plot yield estimation workflow. CPU times were measured on the basis of the 5 validation  
538 datasets and were artificially transformed into an average time per image, while the processing  
539 of foreground trees alone included an additional post-processing script to exclude expert  
540 annotations and network detections outside of the polygons circumscribing the foreground  
541 trees. Otherwise, a marked 0.007 s fluctuation in GPU time was noted. This could undoubtedly  
542 be explained by the different complexities of the assessed neural architectures, e.g. the  
543 YOLOv11 Nano architecture had 2 million parameters, as compared to the YOLOv8 eXtra-  
544 large architecture with 68 million parameters. GPU times for cocoa pod detection would  
545 therefore have an impact on the time to count cocoa pods in a plot yield estimation workflow.  
546 Yet the GPU time seemed to be an interesting criterion that could potentially be used to  
547 complement the network classification. RTDet networks were discarded because their  
548 average GPU times per image were 15-20% higher than those of YOLOv8 and YOLOv11,  
549 even though their average cocoa pod detection capacity was only 2% lower than that of  
550 YOLOv8.

551 The persistent and temporary data storage sizes were easy to evaluate for the different neural  
552 architectures when the algorithmic implementations were perfectly mastered. For example, the  
553 weight file size was constant for a given neural architecture, while ranging from a few kilobytes  
554 to several megabytes depending on the case. Systematic data tiling for training and validation  
555 enabled *a priori* evaluation of the temporary data size, while precise specification of the  
556 treatment processes enabled evaluation of the persistent data size. Here we deliberately have  
557 not elaborated on these aspects since it was not possible to obtain a realistic assessment of  
558 the carbon footprint due to the lack of precision in estimating the CPU and GPU times.

559 However, the trends that emerged from this initial research should help us draw up guidelines  
560 for future studies.

561 It would thus be essential to enrich the datasets with truly new annotations representative of  
562 the different field contexts. This would significantly improve the performance and robustness  
563 of the network, while increasing the amount of data to confirm the results drawn from this study  
564 and enable more meaningful statistical studies on processing times to be conducted. A  
565 calibration pattern, e.g. a blue sphere of known diameter, could be introduced to enable image  
566 calibration, which in turn would enhance the precision on cocoa pod distributions by size class.  
567 The trends noted in this first comparative study will now have to be confirmed in further  
568 research, particularly regarding the predominance of the YOLOv8 Nano architecture. If these  
569 trends are confirmed, the real performance of this latter architecture will have to be assessed  
570 with an embedded system so as to obtain all of the elements needed to define the technical  
571 foundations of an operational solution to help estimate the yield of cocoa tree plots.

572

## 573 Conclusion

574 In this study, the lightest and heaviest architectures of the Faster R-CNN, Detr-ResNet Vision  
575 Transformer, Baidu's Real Time Detection Transformer, YOLOv5, YOLOv8, and YOLOv11,  
576 with identical hyperparameters, were trained with 7,850 annotated pods from 400 low-  
577 resolution images of different cocoa plots. Each network architecture was validated on two  
578 complementary datasets that had not been involved in the training phases. Two validation  
579 datasets were defined in accordance with the recommended data acquisition protocol for early  
580 yield estimation, each corresponding to a cocoa plot. The S42 validation dataset comprised 42  
581 images containing 990 annotated cocoa pods, and the S100 dataset comprised 100 images  
582 containing a total of 2,400 annotated cocoa pods. The validations were based on two  
583 modalities: detection of cocoa pods in the entire set of images, and otherwise in the foreground  
584 tree circumscribed by a specific polygonal line.

585 Unsurprisingly, the most recent networks published as of 2023 gave better results than the  
586 older networks. The top three networks were clearly RTDetr, YOLOv8 and YOLOv11. With  
587 average F1 scores close to 48% for detecting cocoa pods in the background, 88% for those in  
588 the middle, and 94.5% for those in the foreground, all with a maximum dispersion of 2%, the  
589 study confirmed that the YOLOv8, YOLOv11 and RTDetr networks were good candidates for  
590 identifying and counting cocoa pods on foreground trees with sufficient accuracy. However,  
591 with an F1 score of ~77% for the detection of cocoa pods in all images and 90% for foreground  
592 cocoa trees, the YOLOv8 Nano architecture outperformed the other two by almost 3%.

593 Beyond this ranking, our study was designed to assess whether the drop in the performance  
594 of Nano architectures—which, unlike the eXtra-large architectures, prioritized speed to the  
595 detriment of precision—was still acceptable for yield estimates. Against all expectations, the  
596 YOLOv8 and YOLOv11 Nano architectures achieved F1-scores that were a few hundredths  
597 higher than those of the eXtra-large YOLO and RTDetr architectures. Consequently, the eXtra-  
598 large architectures—with 20- to 30-times more parameters to adjust than the Nano  
599 architectures—would probably have required larger datasets to be able to set up truly relevant  
600 network training sessions so as to enable fully objective inter-network comparisons.

601 However, one of the major aspects for cocoa tree plot yield estimation was the performance  
602 sensitivity to data from different networks. Two validation sets did not provide certainty, only  
603 trends. With a sensitivity of less than 0.5%, YOLOv8 Nano took the lead in the top three.

604 With a dispersion of <0.007 s, the average CPU times per image between the different  
605 networks were considered to be similar, i.e. this slight fluctuation clearly had no impact on the  
606 cocoa pod detection workflow. With a dispersion of 0.07 s, the average GPU times per image  
607 reflected the impact of variations in the complexity of the architectures of the different networks  
608 studied. The GPU time was a complementary network classification indicator to choose the  
609 best trained network to implement in the cocoa pod counting workflow. With average GPU  
610 times per image 15-20% higher than those of the YOLOv8 and YOLOv11 networks, the RTDetr  
611 Large and eXtra-large architectures were discarded even though their average cocoa pod

612 detection capacity was just 2% lower than that of the Nano YOLOv8. As the persistent and  
613 temporary data storage sizes were easy to evaluate according to the different neural  
614 architectures, the trends that emerged from our preliminary study should help establish  
615 guidelines for future studies to define a low-carbon footprint tool for early yield estimation in  
616 cocoa plantations.

## 617 Acknowledgements

618 This study was conducted within the framework of the Cocoa4Future (C4F) project, which is  
619 funded by the European DeSIRA Initiative under grant agreement No. FOOD/2019/412-132  
620 and by the French Development Agency. L-CB was funded by Phenome project (ANR-11-  
621 INBS-0012). Juan-Pablo Rojas-Bustos is acknowledged for his participation in earlier  
622 discussions related to this work.

623 The authors would also like to thank David Manley for his proofreading and rewriting work.

624 *To Crunchy, my 13-year-old extra-dwarf rabbit who died peacefully in my arms while I was*  
625 *writing the first version of this article.*

## 626 References

- 627 1. Olofintuyi S. S., Olajubu E. A., Olanike D. An ensemble deep learning approach for  
628 predicting cocoa yield. *Heliyon*. 2023;9(4).
- 629 2. Kongor J. E., Owusu M., Oduro-Yeboah C. Cocoa production in the 2020s: Challenges  
630 and solutions. *CABI Agriculture and Bioscience*. 2024;5(1), 102.
- 631 3. Alain B. K. Economic impact of cocoa culture in Ivory Coast and in Ghana from 1980  
632 to 2015. *J His Arch & Anthropol Sci*. 2024;9(2), 62-67.
- 633 4. Gongal A., Amatya S., Karkee M., Zhang Q., Lewis, K. Sensors and systems for fruit  
634 detection and localization: A review. *Computers and electronics in agriculture*.  
635 2015;116, 8-19.
- 636 5. Girshick R., Donahue J., Darrell T., Malik J. Rich feature hierarchies for accurate object  
637 detection and semantic segmentation. *Proceedings of the IEEE conference on  
638 computer vision and pattern recognition*. 2014;580–587.
- 639 6. Hou L., Wu Q., Sun Q., Yang H., Li P. Fruit recognition based on convolution neural  
640 network. *Natural Computation, Fuzzy Systems and Knowledge Discovery (ICNC-  
641 FSKD)*. 2016; 12th International Conference on, (IEEE), 18-22.
- 642 7. Rahnemoonfar M., Sheppard C. Deep count: fruit counting based on deep simulated  
643 learning. *Sensors*. 2017;17(4), 905.
- 644 8. Xiao F., Wang H., Xu Y., Zhang R. Fruit detection and recognition based on deep  
645 learning for automatic harvesting: An overview and review. *Agronomy*. 2023;13(6),  
646 1625.
- 647 9. Borianne P., Sarron J., Borne F., Faye E. Deep mango cultivars: cultivar detection by  
648 classification method with maximum misidentification rate estimation. *Precision  
649 Agriculture*. 2023;24(4), 1619-1637.
- 650 10. Daviet B., Fournier C., Cabrera-Bosquet L., Simonneau T., Cafier M., Romieu C.  
651 Ripening dynamics revisited: an automated method to track the development of  
652 asynchronous berries on time-lapse images. *Plant Methods*. 2023;19(1), 146.

- 653 11. Sa I., Ge Z., Dayoub F., Upcroft B., Perez T., McCool C. Deepfruits: A fruit detection  
654 system using deep neural networks. *Sensors*. 2016;16(8), 1222.
- 655 12. Borianne P., Borne F., Sarron J., Faye É. Deep Mangoes: from fruit detection to cultivar  
656 identification in colour images of mango trees. *arXiv preprint*. 2019;1909.10939.
- 657 13. Wan S., Goudos, S. Faster R-CNN for multi-class fruit detection using a robotic vision  
658 system. *Computer Networks*. 2020;168, 107036.
- 659 14. Mirhaji H., Soleymani M., Asakereh A., Mehdizadeh S. A. Fruit detection and load  
660 estimation of an orange orchard using the YOLO models through simple approaches  
661 in different imaging and illumination conditions. *Computers and Electronics in  
662 Agriculture*. 2021;191, 106533.
- 663 15. Cheng R. A survey: Comparison between Convolutional Neural Network and YOLO in  
664 image identification. In *Journal of Physics: Conference Series*. 2020; Vol. 1453, No. 1,  
665 p. 012139). IOP Publishing.
- 666 16. Sharma A., Kumar V., Longchamps L. Comparative performance of YOLOv8, YOLOv9,  
667 YOLOv10, YOLOv11 and Faster R-CNN models for detection of multiple weed species.  
668 *Smart Agricultural Technology*. 2024;9, 100648.
- 669 17. Khanam R., Asghar T., Hussain M. Comparative Performance Evaluation of YOLOv5,  
670 YOLOv8, and YOLOv11 for Solar Panel Defect Detection. In *Solar*, MDPI. 2025; Vol.  
671 5, No. 1, p. 6.
- 672 18. Wan S., Goudos S. Faster R-CNN for multi-class fruit detection using a robotic vision  
673 system. *Computer Networks*. 2020;168, 107036.
- 674 19. Raj R., Nagaraj S. S., Ritesh S., Thussha, T. A., Aparanji V. M. Fruit classification  
675 comparison based on cnn and yolo. In *IOP Conference Series: Materials Science and  
676 Engineering*. 2021; Vol. 1187, No. 1, p. 012031.
- 677 20. Jrondi Z., Moussaid A., Hadi M. Y. Exploring End-to-End object detection with  
678 transformers versus YOLOv8 for enhanced citrus fruit detection within trees. *Systems  
679 and Soft Computing*. 2024; 6, 200103.

- 680 21. Mamadou D., Kacoutchy J. A., Ballo A. B., Kouassi B. M. Cocoa pods diseases  
681 detection by MobileNet confluence and classification algorithms. *International Journal*  
682 *of Advanced Computer Science and Applications*. 2023;14(9).
- 683 22. Vera D. B., Oviedo B., Casanova W. C., Zambrano-Vega C. Deep learning-based  
684 computational model for disease identification in cocoa pods (*Theobroma cacao* L.).  
685 arXiv preprint arXiv: 2024;2401.01247.
- 686 23. Ayubi A., Faiz M., Situmorang G. B., Ramadhani K. N., Utama N. P. A Cocoa Ripeness  
687 Detection and Classification Model Based on Improved YOLOv5s. In 2023 10th  
688 International Conference on Advanced Informatics: Concept, Theory and Application  
689 (ICAICTA). 2023;1-6. IEEE.
- 690 24. Sykes J. R., Denby K. J., Franks D. W. Computer vision for plant pathology: A review  
691 with examples from cocoa agriculture. *Applications in Plant Sciences*. 2024;12(2),  
692 e11559.
- 693 25. Ayikpa K.J., Mamadou D., Sodjinou S.G., Ballo A.B., Gouton P., Adou K.J. Detecting  
694 and Extracting Cocoa Pods in the Natural Environment Using Deep Learning Methods.  
695 In International Conference on Digital Technologies and Applications, Springer Nature  
696 Switzerland. 2023;164-174
- 697 26. Lammoglia S. K. D., Borianne P., Théveny F., Cabrera-Bosquet L. Real-time Image  
698 detection of cocoa pods in natural environment using deep learning algorithms. *ICCO*.  
699 2023.
- 700 27. Ren S., He K., Girshick R., Sun J. Faster r-cnn: Towards real-time object detection with  
701 region proposal networks. *Advances in neural information processing systems*. 2025;  
702 28, 91-99
- 703 28. Liu Y., Lu B., Peng J., Zhang Z. Research on the use of YOLOv5 object detection  
704 algorithm in mask wearing recognition. *World Sci. Res. J.* 2020; 6(11), 276-284.
- 705 29. Reis D., Kupec J., Hong J., Daoudi A. Real-time flying object detection with YOLOv8.  
706 arXiv preprint arXiv. 2023;:2305.09972.

- 707 30. Khanam R., Hussain M. YOLOv11: An overview of the key architectural enhancements.  
708 arXiv 2024. arXiv preprint arXiv:2024;2410.17725.
- 709 31. Carion N., Massa F., Synnaeve G., Usunier N., Kirillov A., Zagoruyko, S. End-to-end  
710 object detection with transformers. In European conference on computer vision  
711 2020:(pp. 213-229). Cham: Springer International Publishing.
- 712 32. Zhao Y., Lv W., Xu S., Wei J., Wang G., Dang Q., ... Chen J. Detsr beat yolos on real-  
713 time object detection. In Proceedings of the IEEE/CVF conference on computer vision  
714 and pattern recognition. 2024;16965-16974.
- 715 33. Jegham N., Koh C.Y., Abdelatti M., Hendawi A. Evaluating the evolution of yolo (you  
716 only look once) models: A comprehensive benchmark study of yolo11 and its  
717 predecessors. arXiv preprint arXiv. 2024;2411.00201.
- 718 34. Arkin E., Yadikar N., Xu X., Aysa A., Ubul, K. A survey: object detection methods from  
719 CNN to transformer. Multimedia Tools and Applications. 2023;82(14), 21353-21383.
- 720 35. He K., Zhang X., Ren S., Sun, J. Deep residual learning for image recognition. In  
721 Proceedings of the IEEE conference on computer vision and pattern recognition.  
722 2016;770-778.
- 723 36. Benaben D., Berthoud F., Guennebaud G., Ligozat A.L., Valcke S. Estimation de  
724 l'empreinte carbone d'une heure de calcul sur un cœur CPU ou sur un GPU (Doctoral  
725 dissertation, Labos 1point5). 2024.
- 726 37. Oddoye E.O., Agyente-Badu C.K., Gyedu-Akoto E. Cocoa and its by-products:  
727 Identification and utilization. Chocolate in health and nutrition. 2013;23-37.
- 728 38. Sabas B.Y.S., Danmo K.G., Madeleine K.A.T., Bogaert J. Cocoa production and forest  
729 dynamics in ivory coast from 1985 to 2019. Land, 2020;9(12).
- 730 39. Dago M. R., Zo-Bi I. C., Konan I. K., Kouassi A. K., Guei S., Jagoret P., Hérault B. What  
731 motivates West African cocoa farmers to value trees? Taking the 4 W approach to the  
732 heart of the field. People and Nature. 2025;7(1), 215-230.
- 733 40. Lammoglia S.K.D., Cabrera-Bosquet L., Rojas-Bustos J.P. 2022.  
734 <https://www.zooniverse.org/projects/phenoarch/cocoa-fruit-counting>

- 735 41. Simpson R., Page K. R., De Roure D. Zooniverse: observing the world's largest citizen  
736 science platform. In Proceedings of the 23rd international conference on world wide  
737 web. 2014;1049-1054.
- 738 42. Schneider C.A., Rasband W.S., Eliceiri K.W. NIH Image to ImageJ: 25 years of image  
739 analysis. Nat Methods. 2012;9, 671–675.
- 740 43. Rahman M.A., Wang Y. Optimizing intersection-over-union in deep neural networks for  
741 image segmentation. In International symposium on visual computing. 2016;234-244.  
742 Cham: Springer International Publishing.
- 743 44. Yacouby R., Axman D. Probabilistic extension of precision, recall, and f1 score for more  
744 thorough evaluation of classification models. In Proceedings of the first workshop on  
745 evaluation and comparison of NLP systems. 2020;79-91..
- 746 45. Borianne P., Théveny F., Bertrand B., Villain L., Faye É., Sarron J., Viennois G., Borne  
747 F., Jaeger J. L'IA au service de l'agriculture : au cœur de l'expérience Agro'Deep. 2021.  
748 (hal-05019239)
- 749 46. Kurtzer G.M., Sochat V., Bauer M.W. Singularity: Scientific containers for mobility of  
750 compute. PloS one, 2017;12(5), e0177459.
- 751 47. Bailly A., Blanc C., Francis É., Guillotin T., Jamal F., Wakim B., Roy P. Effects of  
752 dataset size and interactions on the prediction performance of logistic regression and  
753 deep learning models. Computer Methods and Programs in Biomedicine. 2022;213,  
754 106504.
- 755 48. Vasconez J.P., Delpiano J., Vougioukas S., Cheein F.A. Comparison of convolutional  
756 neural networks in fruit detection and counting: A comprehensive evaluation.  
757 Computers and Electronics in Agriculture. 2020;173, 105348.
- 758 49. Dawson H.L., Dubrule O., John C. M. Impact of dataset size and convolutional neural  
759 network architecture on transfer learning for carbonate rock classification. Computers  
760 & Geosciences. 2023;171, 105284.