# Detecting and Counting Pistachios based on Deep Learning

**Mohammad Rahimzadeh**
School of Computer Engineering
Iran University of Science and Technology, Iran
mr7495@yahoo.com
ORCID: 0000-0002-8550-8967
Corresponding author

**Abolfazl Attar**
Department of Electrical Engineering
Sharif University of Technology, Iran
attar.abolfazl@ee.sharif.edu
ORCID: 0000-0001-6727-432X

May 5, 2021

## Abstract

Pistachios are nutritious nuts that are sorted based on the shape of their shell into two categories: Open-mouth and Closed-mouth. The open-mouth pistachios are higher in price, value, and demand than the closed-mouth pistachios. Because of these differences, it is considerable for production companies to precisely count the number of each kind. This paper aims to propose a new system for counting the different types of pistachios with computer vision. We have introduced and shared a new dataset of pistachios, including six videos with a total length of 167 seconds and 3927 labeled pistachios. Unlike many other works, our model counts pistachios in videos, not images. Counting objects in videos need assigning each object between the video frames so that each object be counted once. The main two challenges in our work are the existence of pistachios' occlusion and deformation of pistachios in different frames because open-mouth pistachios that move and roll on the transportation line may appear as closed-mouth in some frames and open-mouth in other frames. Our novel model first is trained on the RetinaNet object detector network using our dataset to detect different types of pistachios in video frames. After gathering the detections, we apply them to a new counter algorithm based on a new tracker to assign pistachios in consecutive frames with high accuracy. Our model is able to assign pistachios that turn and change their appearance (e.g., open-mouth pistachios that look closed-mouth) to each other so does not count them incorrectly. Our algorithm performs very fast and achieves good counting results. The computed accuracy of our algorithm on six videos (9486 frames) is 94.75%.

***Keywords*** Deep learning · Convolutional Neural Network · Pistachio Counting · Multi-Object Counting · Object Detection · Motile-Object Counting

## 1 Introduction

Nowadays, automation in the industry plays a significant role in increasing efficiency and saving resources. One of the industries that need more development in automation than other industries is the agricultural industry and related fields. Proper packaging of agricultural products will increase profitability and reduce crop losses. On the other hand, crop quality categorization depends on human resources, which causes time-consuming and rising costs, and most importantly, does not have the necessary quality compared to machines.

Pistachio is one of the crops that need human resources to classify and count so that the quality of the crop can be evaluated in terms of its open or closed shell. Pistachios are mostly sorted based on their shell's shape to open-mouth and closed-mouth, and these two kinds differ in price and value.

Pistachios are used as nuts, and in the food industry [39]. Pistachio kernels are rich in unsaturated fatty acids, fiber, carbohydrates, proteins, and various vitamins that are very useful for the human diet [23, 14]. Adequate consumption

of pistachio kernels reduces the risk of heart disease and has a good effect on blood pressure in people who do not have diabetes, and prevents some cancers[9, 14, 38]. Pistachio is one of Middle Eastern countries' main agricultural products, especially Iran, [7]. The largest producers of this product in the world are Iran, the USA, and Turkey, respectively [7].

There are many types of pistachios, depending on the type and place of growth; they have different sizes, colors, and flavors [30]. Depending on the shape of the pistachio, it can be divided into three general categories: round, long, and jumbo [30]. Long pistachios have a narrower split than the other two, and the round and jumbo type have a much clearer split than the semi-closed one [30]. Fig. 1 shows a summary of pistachios' different types.
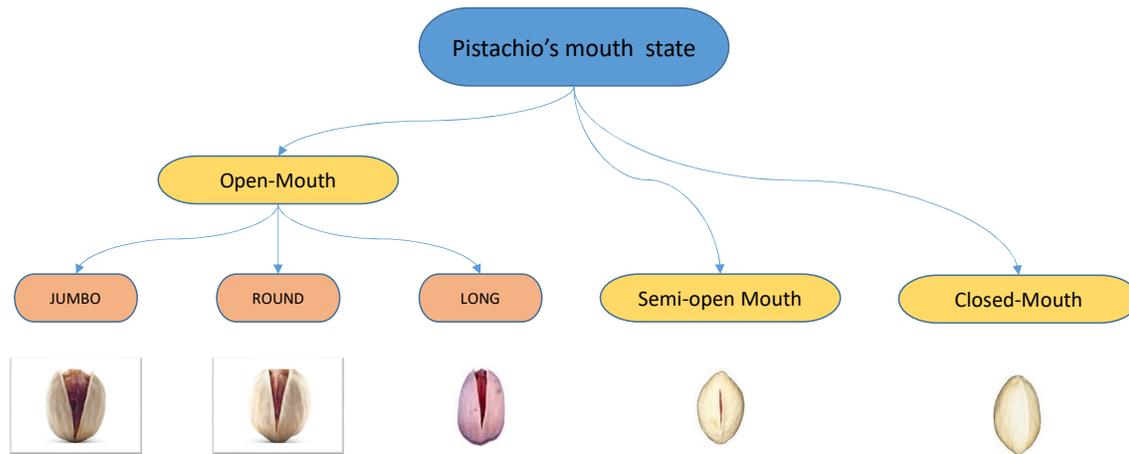


Figure 1: Pistachios Assortment

The average weight of each pistachio's shell is about 0.57 grams [15], which is about 1750 per kilogram. As a result, according to the statistics provided, counting them is a very time-consuming and tedious task that artificial intelligence can easily do.

Detecting and counting the pistachios can be used for proper packaging and crop quality assessment. Another advantage of this can be the estimation of the amount of the crops in the coming years and the breeding of pistachio trees to increase the quality of the crop. As there is a significant difference in price and demand between the open-mouth and closed-mouth pistachios, the factories related to pistachio production or packaging need to know precisely how much of these two kinds exist in every package. Counting these two kinds of pistachios can also help separate them to increase the exporting packages' quality. Counting the pistachios by human resources is very time-consuming and practically uneconomical, so machine vision can play a significant role in this regard.

Another application of these procedures is that closed pistachios are used by the mechanical opening method [3] to be returned to the consumption cycle. In this method, it is first necessary to identify the closed pistachios, which results in reduced losses and increased crop yields [3].

One of the new methods for detecting, counting, and classifying pistachios is machine vision. In recent years, machine vision has been used for many tasks to automate and replace machines with humans, which have yielded excellent results [32]. These applications exist in the fields of medicine [18], medical image diagnosis [26, 28], self-driving [10], security [2], and agriculture [29, 22, 24], and so on.

One of the principles of using robots and remote control and sense is the use of machine vision. Therefore, improving the accuracy and precision of the system is one of the essential principles. In machine vision, various methods such as thermal cameras, sensors, microscopes, and common cameras have been used for imaging space around it. However, the main issue in machine vision is the choice of the data analysis method.

Currently, one of the most attractive and accurate methods of machine vision is deep learning, which has been created a revolution in artificial intelligence [17]. One of the most important advantages of deep convolutional neural networks is the comprehensive and flexible recognition of different objects. [16].

Using the deep convolutional neural networks, we can identify and count pistachios. The appearance of the pistachio, the angle of the camera or robot also plays a vital role in correctly determining the open and closed pistachios.

Our work is the first to investigate the detecting, tracking, and counting performance of pistachios on transportation lines, even in high density and occlusion situations. Another critical challenge is to count the open-mouth and closed-mouth

pistachios correctly because the open-mouth pistachios can show themselves as closed-mouth pistachio when moving and spinning on the transportation line. Based on this situation, many tracking algorithms will fail because the class type of pistachio changes between successive frames. Still, our tracking and counting model has solved this problem and achieves high accuracy.

In this paper, at first, we will propose our dataset, which we call Pesteh-Set. At the next stage, we will describe the detection phase. We have implemented the RetinaNet [20] object detector for detecting the Pistachios in video frames. We have separated the dataset into five-folds and allocated 20 percent of the dataset for testing and the rest for the training. After the detection phase, we present the model for counting the open-mouth and closed-mouth pistachios. This algorithm runs very fast with high accuracy. The general schematic of our work is presented in Fig. 2.
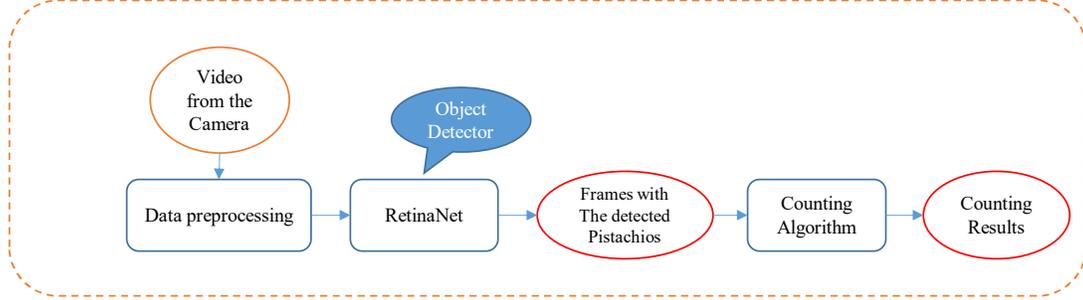


Figure 2: General Schematic of our proposed model



Figure 3: The General View of how Pesteh-Set was recoreded and our proposed system for counting the pistachios

One of the closest studies in our research is fruit detection and counting [11]. This is done using a variety of tools such as a B/W camera [37], a color camera [5], a thermal camera [35] and a spectral camera [31]. Due to the type of data we have, color is one of the most important features, B/W camera is not suitable [37]. On the other hand, because the spectral camera has a time delay, this method cannot be suitable [31]. Due to its sensitivity to size and lack of detection of split pistachios, the thermal camera is not suitable for analyzing our data [35]. As a result, the color camera is more suitable than other tools. Another advantage of using color cameras is their abundance, especially in mobile phones, which can be used for remote monitoring and control.

There are other methods, such as sensors, that can be used for counting, but since one of the challenges ahead is split and non-split pistachios counting, however, this method is not suitable for our data [8].

One of the most widely used classifications and detection methods used is the K-means clustering method, which is performed unsurprisingly. In [36], K-means clustering detects green apples using thermal and color cameras.

One of the most basic methods of Supervised classification is the Bayesian classifier, which has been used in [34] to identify oranges and has yielded relatively good results from previous research. Other used methods include KNN clustering [21].

Artificial neural networks have a special place in machine vision and object detection. In the meantime, deep convolutional neural networks have shown excellent results for images and videos, too. Therefore, image detection systems move to increase the accuracy and quality of deep learning.

In [24], they classify different fruits using an innovative deep neural network. In [29], a deep neural network was developed to detect and count the number of tomatoes per plant. In this study, due to the lack of enough data, the data were generated by simulating a green and brown environment in which the tomatoes were simulated with red circles. This can take the results away from the real world.

In [22], the environment is photographed using a monocular camera, and visible fruits are detected and tracked on the tree. This detection is made by training a fully Convolutional Network, then using image processing methods to track the fruits and then counting them.

[27] is one of the papers that has done very well in detecting and tracking motile objects. This paper also introduced a method for improving motile-objects detection. In this article, by using RetinaNet [20], and other introduced methods, they have detected motile sperms in the video frames. Finally, by implementing a new tracker called modified CSR-DCF, they have tracked and analyzed the sperms attributes, such as their number and motility characteristics.

The rest of the paper is organized as follows: In section 2, we describe our dataset and the detection and counting models. In section 3, the results of our work are presented. In section 4, we discuss the obtained results, and in section 5, the paper is concluded.

## 2  Materials and methods

### 2.1  Pesteh-Set

Pistachio is known as Pesteh in Iran; that is why we called our dataset Pesteh-Set. Pesteh-Set [1] is made of two parts. The first part includes 423 images with ground truth. We sorted the pistachios into two classes: Open-mouth and closed-mouth. The images' ground truth consists of the bounding boxes of the two classes of pistachios in the images. There are between 1 to 27 pistachios in each image and 3927 pistachios totally. The second part includes six videos (9486 frames) that were used for the counting phase. These six videos include 561 motile pistachios and more than 350,000 single pistachios (sum of pistachios in all frames).

The videos of the dataset have been recorded by a cell-phone camera with $1920 \times 1080$ pixels resolution. Five of these videos are recorded with 60 frames per second(fps) frame rate, and one other is recorded with 30 fps frame rate. The cell-phone was perched on the wall above the line that was transporting the pistachios. This line was designed somehow that the pistachios could roll on it. The reason the pistachios rolling is so important is that the open-mouth pistachios could appear on their backside where they look like closed-mouth pistachios, but the rolling cause them to show their open-mouth side when rolling. Fig. 3 presents a view of how the dataset was recorded and the general schematic of our proposed system for remote counting the pistachios.

We have selected some frames of the videos and labeled them with a self-developed program using OpenCV library [25] on python language. The pistachios are categorized into two classes: open-mouth pistachios and closed-mouth pistachios. Some of the images of this dataset are presented in Fig. 4. The self-developed program for labeling the images along with all the data has been shared so other researchers could use them to make the Pistachio-Dataset larger. Table. 1 presents the details of Pesteh-Set.

In Table. 1, the reason that the number of pistachios in the videos is less than the images is that the number of pistachios in the videos denotes the number of mobile pistachios. It means, from the moment one pistachio enters the video in a frame until it exits the video in the later frames, it would be counted as one pistachio.

To report the number of pistachios in the images, we counted the number of pistachios in each image. It is noteworthy that we selected non-consecutive frames of different videos for labeling, so there would not be a similarity between them (for training the models). Besides, we tried to choose the frames somehow that we have an almost equal number of each class, and our dataset become balanced for training.

---

[1]<This dataset is shared in https://github.com/mr7495/Pesteh-Set>

Figure 4: Some of the images in Pesteh-Set

Table 1: This table shows the distribution of Pistachios in Pesteh-Set

| | Number of Open-mouth Pistachios | Number of Closed-mouth Pistachios | Number of All the Pistachios |
|---|---|---|---|
| Video 1 | 50 | 20 | 70 |
| Video 2 | 60 | 20 | 80 |
| Video 3 | 70 | 20 | 90 |
| Video 4 | 90 | 20 | 110 |
| Video 5 | 100 | 20 | 120 |
| Video 6 | 39 | 52 | 91 |
| All of the Videos | 409 | 152 | 561 |
| All the 423 labeled images | 1993 | 1934 | 3927 |

## 2.2 Detection

### 2.2.1 RetinaNet

RetinaNet [20] is a deep fully convolutional neural network that is utilized for object detection. The architecture of RetinaNet is depicted in Fig. 5. RetinanNet is made of three main parts. The first part, which is the feature extractor, is build up from a backbone model and the feature pyramid network (FPN) [19]. Famous convolutional networks like ResNet [12], DenseNet [13], and VGG [33] are mostly used as the backbone of RetinaNet. The FPN takes the multi-dimensional features that are extracted from the backbone network as the input to build a multi-scale feature pyramid from the input image [19]. The usage of the FPN on top of the backbone considerably improves object detection accuracy because it makes the model able to detect multi-scale objects.
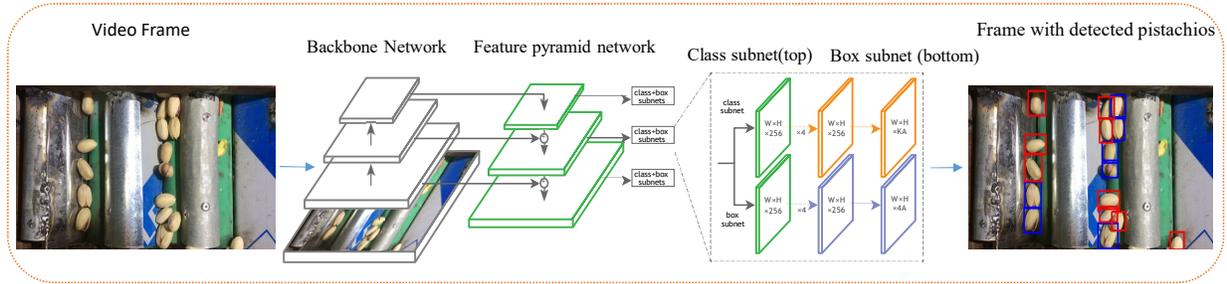
5

Figure 5: RetinaNet Architecture

The second part of the RetinaNet is the Classifier, which has the role in predicting the possibility of the presence of each of the classes at each spatial location for each of the anchor boxes. The third part is the box regression that regresses each of the anchor boxes to the nearest ground truth object boxes [20]. Another novelty presented in the RetinaNet is using the focal loss [20] as the loss function. The focal loss adds a modulating factor to the cross-entropy loss function, and by doing so, it focuses on the hard examples while training, and as the loss of hard examples is higher than the easy examples, it improves the learning process and accuracy.

### 2.2.2 Training

We have separated the Pesteh-Set into five folds for training, which in each fold, 20 percent of the dataset was allocated for testing, and the rest for training. The original size of images in our dataset is 1070x600 pixels, but they were preprocessed and then resized to 1333×747 pixels for training and testing the object detectors.

We used RetinaNet [20], as the object detector. We trained and validated RetinaNet on 3 different backbones: ResNet50 [12], ResNet152 [12], and VGG16 [33]. VGG16 [33] was introduced by Oxford researchers and in 2014 and is made of 13 layers for the feature extraction part and represents a feature map with 512 channels. ResNet [12] series was introduced after the VGG and introduced residual layers that improved deep convolution networks. These layers help the networks avoid losing data when the features being corrupted while passing through layers by applying skip connections from the previous layers to the next layers. Since ResNet, most of the next model also used and inspired from residual layers. ResNet comes with several models like ResNet50, ResNet101, and ResNet152; the difference between them is having more layers, and all of them represent a feature map with 2048 channels. ResNet152 achieved higher results on the ImageNet dataset [6]. The details of RetinaNet on these backbone models are described in Table. 2.

Table 2: The number of parameters and training, and inference time of RetinaNet on different backbones are listed in this table.

| Backbone | Number of Parameters | Training Time per Epoch | Inference Time per Image with Size (1333,747,3) |
|---|---|---|---|
| VGG16 | 23,428,470 | 318 s | 17.44 ms |
| ResNet50 | 36,403,702 | 390 s | 17.80 ms |
| ResNet152 | 71,137,782 | 615 s | 20.01 ms |

Transfer learning from the ImageNet [6] pre-trained weights was utilized on the backbone model at the beginning of the training to speed up the network convergence. We also used data augmentation methods like rotation, translation, shearing, horizontal and vertical flipping, and rescaling to improve the learning efficiency and stop the network from overfitting.

We applied Adam optimizer with a 1e-5 initial learning rate and an automatic reduction function to reduce the learning rate when training loss cannot be improved for five epochs. The focal loss was used for the classification subnet, and the Smooth L1 loss function was used for the regression subnet.

In classification problems, it is most usual to set batch size between 32 to 256, so the model can analyze several images in each step of training and learn the features better. In object detection tasks, if several objects are available in each image, even if the batch size is equivalent to 1, the model analyzes several objects in each step of training. That means

the actual batch size is equal to the number of objects in that image. Based on this fact, as there are several pistachios in each image, we set the batch size equal to 1 to reduce computational costs.

The applied training parameters are listed in the Table. 3 and the details of each fold are present in Table. 4.

Table 3: This table shows all the parameters and methods we used in training

| Training Parameters | Value |
|---|---|
| Learning Rate | 1e-5 (With automatic reduction based on Loss value) |
| Batch Size | 1 |
| Optimizer | Adam |
| Loss Function | Focal Loss for the classification subnet<br>Smooth L1 for the regression subnet |
| Steps | 1017 |
| Horizontal/Vertical flipping | Yes (50%) |
| Translation Range | -0.1 - 0.1 |
| Rotation Range | 0 - 360 degree |
| Shear Range | -0.1 - 0.1 |
| Scaling Range | -0.1 - 0.1 |

Table 4: This table presents the details of our train and test sets in each fold

| Fold Number | Training Images | Test Images | Open-mouth Pistachios in Training Set | Closed-mouth Pistachios in Training Set | Open-mouth Pistachios in Testing Set | Closed-mouth Pistachios in Testing Set |
|---|---|---|---|---|---|---|
| Fold 1 | 339 | 84 | 1600 | 1550 | 393 | 384 |
| Fold 2 | 339 | 84 | 1610 | 1572 | 383 | 362 |
| Fold 3 | 339 | 84 | 1553 | 1506 | 440 | 428 |
| Fold 4 | 339 | 84 | 1641 | 1575 | 352 | 359 |
| Fold 5 | 336 | 87 | 1568 | 1533 | 425 | 401 |

## 2.3   Counting

The second and main phase of our work was counting the number of open-mouth and closed-mouth pistachios in the videos. There were several challenges in this phase. The first challenge was that we wanted to develop a counting method that could be performed very fast on the CPU. Some of the other ideas may need a GPU; otherwise, the process would become highly time-consuming. However, our method works very much fast on the CPU, even faster than many other methods on GPU. As the object detection part needs a GPU for running in real-time mode, by applying a counting model on the CPU, the process of detection and counting can be executed simultaneously and much faster with no need to providing more expensive hardware.

The second challenge was that some of the open-mouth pistachios could show themselves as closed-mouth pistachio in some consecutive frames and then reveal their open part only a few frames. Moreover, some open-mouth pistachios rolling on the transportation line could show their open part several times and then appear like closed-mouth pistachios like Fig. 7. We had to develop our algorithm somehow to prevent failing because of these challenges.

Another challenge is developing the counting method to prevent failure because of false detections or not-detected pistachios that may be affected by the pistachios' occlusion.

For counting the pistachios, we first generate the frames from the taken video and then use the trained network for getting the bounding boxes of the pistachios in the frames. After this process, we would have a list of bounding boxes from the video's frames.

Two thresholds have been designed to improve the counting accuracy: the initial threshold and the end threshold. The initial threshold is set to detect the newly inserted pistachios, and the end threshold is to reject adding the pistachios to the tracks list. These two methods are explained in the next paragraphs.

The algorithm first runs a function to set the initial threshold. In this function, the algorithm begins to assign the pistachios between each of two consecutive frames based on their distance without considering the class of pistachios. The minimum acceptable distance to assign the pistachios has been set to 20 pixels. The number of 20 pixels is taken out of our experience to prevent false assignments. After the assigning, the function adds the not-assigned pistachios that the height of the mid-point of their bounding box be less than 200 to a list (the height of the images is 600 pixels). These added pistachios are candidates as new inputs. After adding all the eligible pistachios, the function measures the average of the list, and it would be set as the initial threshold. This process performs to measure the area that most of the pistachios will enter the frame. The pistachios in different videos can enter the video frames differently and also may have various speeds, so this function will set the initial threshold separately for each video to improve the counting performance.

As the height of our image is 600 pixels, we set the end threshold equal to 500. We call the area between the top of an image (y=0) to the initial threshold (indicated by the counting model), the entering area, and the area between the end threshold (y=500) to the end of the image (y=600) the exiting area. The entering area is the main region for adding new incoming pistachios, but the exiting area includes pistachios that have been counted before, so we ignore them.

In the next level, the algorithm uses the assigned pistachios of each of the two consecutive frames that were computed in the last step. Our counter algorithm role is to count the number of all pistachios and count the number of open-mouth and closed-mouth pistachios. The algorithm uses the initial threshold and the end threshold, which are computed in the last step. Toward solving the main challenge, which is that many of the pistachios may show their open part in some frames and the closed parts in the other frames, we have to track them by assigning them from frame to frame. We decided to track them from when they enter the entering area until they enter the exiting area. By doing this, we can know if one pistachio is open-mouth or closed-mouth, and it also prevents the algorithm from counting extra open-mouth pistachios.

Our algorithm analyzes each of the assigned pistachios in the two consecutive frames for all the frames. If the pistachio in the last frame was in the existing area or this pistachio in the current frame is in the exiting area, this pistachio would also be rejected to be added to the track list. Otherwise, the pistachio in the current frame would be added to the track list that its assigned pistachio in the previous frame exists in it.

After adding all the assigned pistachio in the current frame to the track lists that they belong to, the algorithm investigates the pistachios in the current frame that have not been assigned to any other pistachios. If these pistachios are located in the exiting area, they would be rejected for adding to any track list. If they be in the entering area and the number of all the pistachios in the current frame be greater than the number of pistachios in the last frame, these pistachios would be considered as new inputs and would start their own track list. The reason the pistachios in the current frame must be higher than the pistachios in the last frame is that in most cases, if a new pistachio enters the frame, the number of all pistachios should increase, but you may think that this situation may not always happen. It is true, but it also equalizes the conditions that some pistachios in the entering area will not be considered new inputs several times mistakenly. The unassigned pistachios in the current frame that can not be chosen as new inputs would be added to the Lost-Pistachios list.

The Lost-Pistachios list is created to assign the pistachios that could not be assigned to the last frame pistachios (maybe because the pistachios in the last frame are not detected) to the pistachios in the 2 to 6 previous frames. If the assignment is successful, the newly assigned pistachio will be added to the track list, and if not, they will be rejected.

Finally, after repeating this procedure for all the consecutive frames, we would have a list of tracked pistachios. If there be an open-mouth pistachio in a track, the whole track will be considered open-mouth. Therefore we could count the open-mouth and closed-mouth pistachios. The flow chart of the proposed counting algorithm is presented in Fig. 6.
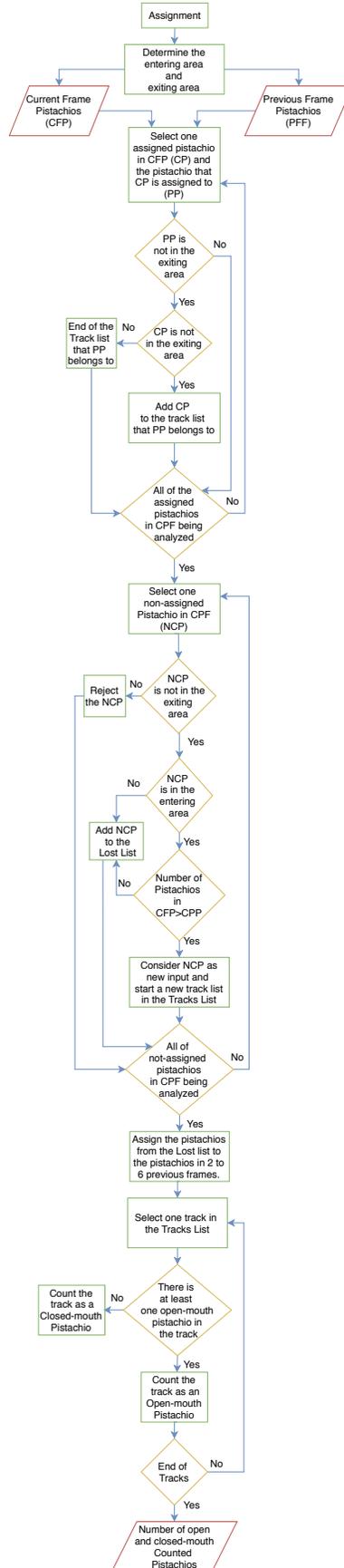
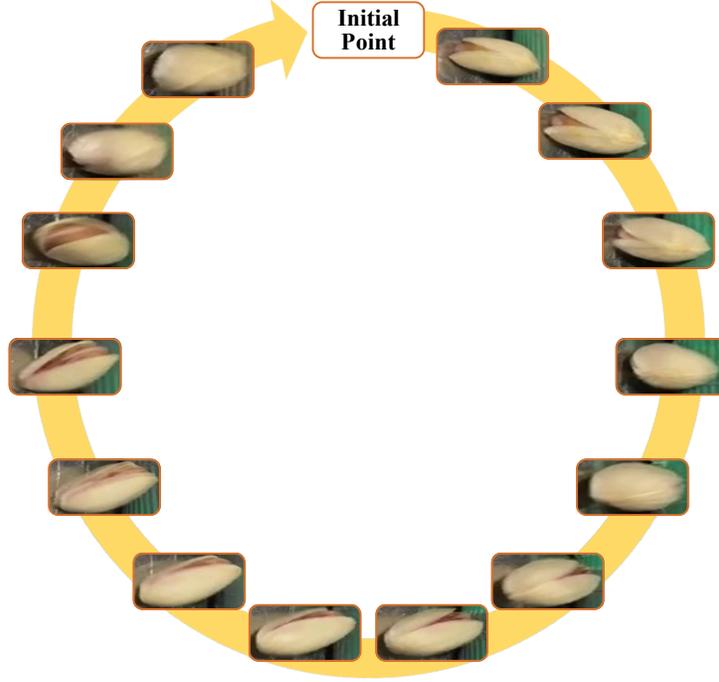Figure 6: The flow chart of the proposed counting algorithm

Figure 7: In this figure, you can observe that a pistachio can be presented as open-mouth and closed-mouth several times while moving.

## 3 Results

### 3.1 Detection Results

We trained Retinanet with three different backbones: ResNet50, ResNet152, and VGG16 based on the explained parameters in Table. 4 for 50 epochs. Data augmentation methods like rotation, translation, shearing, horizontal and vertical flipping, and rescaling were also applied to improve the training and prevent overfitting.

The system we used in this paper was provided by Google Colaboratory Notebooks, which allocated a Tesla P100 GPU, 2.00GHz Intel Xeon CPU, and 25GB RAM on Linux to us. For utilizing RetinaNet we used the written codes by Fizyr which implemented RetinaNet with Keras library [4] on Tensorflow backend [1]. The metrics we used for evaluating RetinaNet in the detection phase are Recall, Precision, F1 score, Accuracy, Average Precision(AP), and Mean Average Precision(mAP). AP is defined as:

$$AP = \frac{\sum_{i=1}^{D} \{Precision(i) \times Recall(i)\}}{Number\ of\ annotations} \tag{1}$$

In Eq. 1, D is the number of detected pistachios that are sorted by scores. Average Precision will be calculated for each class separately. The mAP metric is the mean of Average Precision between classes. Fig. 8 presents some of the images with the detected pistachios.

We considered the detected boxes with Intersection over Union (IOU) more than 0.5 as true positives and the others as false positives to evaluate the detection results. The detection results of RetinaNet on the three backbones are presented in Tables. 6 and 7.

Table. 7. shows that the performance of ResNet backbones is better than VGG16, as expected. The average mAP between five-folds for VGG16, ResNet152, and ResNet50 is 91.23%; 91.69%, and 91.87% respectively. The performance of ResNet50 and ResNet152 is nearly close. Understandably, as there are only two classes in this dataset, and pistachios attributes are distinguishable from the background, there is no need to use bigger models. Using bigger models in our dataset may lead to overfitting and does not increase the output mAP. The VGG16 model is much weaker in image classification tasks than the ResNet models on the ImageNet dataset. Still, RetinaNet on the VGG16 backbone model obtained good detection results, which shows the high effect of feature pyramid network and focal loss.
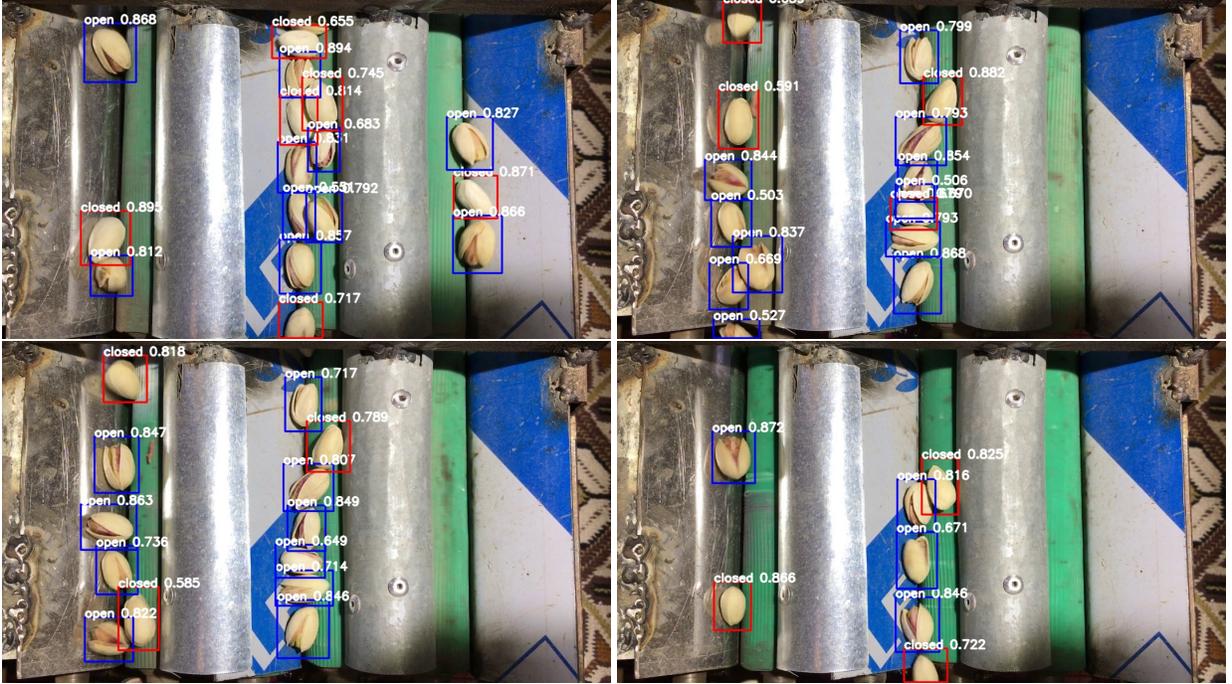
Figure 8: These images are the output of RetinaNet. The red boxes are the closed-mouth pistachios, and the blue boxes belong to the open-mouth pistachios. The number beside the open or closed is the detection score value.

## 3.2 Counting Results

Six different videos with 167 seconds length and 9486 frames were selected for evaluating the counting algorithm. We tested our counting algorithm based on the detections gathered from the trained networks on different backbones. The results and the overall accuracy for all the videos are present in Table. 8 and Table. 5 also expresses the running time of the counting model.

In Table. 8, the accuracy metric is considered for evaluating the tracking algorithm, which is defined as

$$Accuracy = \frac{TP}{TP + FN + FP} \tag{2}$$

In equation 2, TP is the number of the correct-counted pistachios, FN is the number of not-counted pistachios, and FP is the number of extra miscounted pistachios.

Table 5: In this table, the processing time of our counter algorithm after getting the detections from RetinaNet is reported for each video.

| Video | Number of Frames | Time (S) |
|---|---|---|
| Video1 | 984 | 3.00 |
| Video2 | 1665 | 3.39 |
| Video3 | 1833 | 5.11 |
| Video4 | 2227 | 5.75 |
| Video5 | 2171 | 4.19 |
| Video6 | 606 | 0.83 |

11

Table 6: The detection results of RetinaNet on different backbones for each fold are reported in this table.

| | | Class closed-pistachio | | | Class open-pistachio | | |
|---|---|---|---|---|---|---|---|
| | | AP | F1 score | Recall | AP | F1 score | Recall |
| Fold1 | ResNet50 | 0.9072 | 0.9128 | 0.9270 | 0.9467 | 0.9325 | 0.9669 |
| | ResNet152 | 0.9037 | 0.9238 | 0.9166 | 0.9686 | 0.9302 | 0.9847 |
| | VGG16 | 0.8821 | 0.9067 | 0.8984 | 0.976 | 0.9397 | 0.9923 |
| Fold2 | ResNet50 | 0.9135 | 0.9213 | 0.9226 | 0.9172 | 0.9359 | 0.9347 |
| | ResNet152 | 0.9010 | 0.9105 | 0.9143 | 0.9151 | 0.9286 | 0.9347 |
| | VGG16 | 0.8841 | 0.8856 | 0.9198 | 0.9157 | 0.8988 | 0.9399 |
| Fold3 | ResNet50 | 0.9402 | 0.9232 | 0.9556 | 0.8832 | 0.9130 | 0.9068 |
| | ResNet152 | 0.9096 | 0.9161 | 0.9322 | 0.8981 | 0.9171 | 0.9181 |
| | VGG16 | 0.8919 | 0.9138 | 0.9042 | 0.8922 | 0.9032 | 0.9227 |
| Fold4 | ResNet50 | 0.9183 | 0.9115 | 0.9331 | 0.9301 | 0.9258 | 0.9403 |
| | ResNet152 | 0.9347 | 0.9243 | 0.9526 | 0.9449 | 0.9318 | 0.9517 |
| | VGG16 | 0.9186 | 0.9110 | 0.9415 | 0.9303 | 0.9194 | 0.9403 |
| Fold5 | ResNet50 | 0.9122 | 0.9106 | 0.9276 | 0.9179 | 0.9175 | 0.9294 |
| | ResNet152 | 0.8830 | 0.8996 | 0.9052 | 0.9100 | 0.9160 | 0.9247 |
| | VGG16 | 0.8922 | 0.8907 | 0.9152 | 0.8979 | 0.9153 | 0.9035 |

Table 7: This table contains the RetinaNet evaluation data for all the classes.

| | Backbone Network | TP | FP | FN | Recall | Precision | F1 score | mAP | Accuracy |
|---|---|---|---|---|---|---|---|---|---|
| Fold1 | ResNet50 | 736 | 82 | 41 | 0.9472 | 0.8997 | 0.9228 | 0.9270 | 0.8568 |
| | ResNet152 | 739 | 78 | 38 | **0.9510** | **0.9045** | **0.9272** | **0.9361** | **0.8643** |
| | VGG16 | 735 | 79 | 42 | 0.9459 | 0.9029 | 0.9239 | 0.9291 | 0.8586 |
| Fold2 | ResNet50 | 692 | 53 | 53 | 0.9288 | **0.9288** | **0.9288** | **0.9153** | **0.8671** |
| | ResNet152 | 689 | 64 | 56 | 0.9248 | 0.9150 | 0.9198 | 0.90812 | 0.8516 |
| | VGG16 | 693 | 115 | 52 | **0.9302** | 0.8576 | 0.8924 | 0.8999 | 0.8058 |
| Fold3 | ResNet50 | 808 | 84 | 60 | **0.9308** | 0.9058 | **0.9181** | **0.9117** | **0.8487** |
| | ResNet152 | 803 | 81 | 65 | 0.9251 | **0.9083** | 0.9166 | 0.9038 | 0.8461 |
| | VGG16 | 793 | 85 | 75 | 0.9135 | 0.9031 | 0.9083 | 0.8921 | 0.8321 |
| Fold4 | ResNet50 | 666 | 73 | 45 | 0.9367 | 0.9012 | 0.9186 | 0.9242 | 0.8494 |
| | ResNet152 | 677 | 71 | 34 | **0.9521** | 0.9050 | **0.9280** | **0.9398** | **0.8657** |
| | VGG16 | 669 | 82 | 42 | 0.9409 | 0.8908 | 0.9151 | 0.9244 | 0.8436 |
| Fold5 | ResNet50 | 767 | 85 | 59 | **0.9285** | 0.9002 | **0.9141** | **0.9151** | **0.8419** |
| | ResNet152 | 756 | 83 | 70 | 0.9152 | **0.9010** | 0.9081 | 0.8965 | 0.8316 |
| | VGG16 | 751 | 86 | 75 | 0.9092 | 0.8972 | 0.9031 | 0.8950 | 0.8234 |
| Average | ResNet50 | **733.8** | **75.4** | **51.6** | **0.9344** | **0.9071** | **0.9205** | **0.9187** | **0.8528** |
| | ResNet152 | 732.8 | **75.4** | 52.6 | 0.9336 | 0.9068 | 0.9199 | 0.9169 | 0.8519 |
| | VGG16 | 728.2 | 89.4 | 57.2 | 0.9279 | 0.8903 | 0.9086 | 0.9123 | 0.8332 |

The information in Table. 8, have been calculated based on the trained models in our dataset's first fold. As ResNet152 got better results in Fold1 (Table. 7), we can see that the counting results are better for this model in Fold1. The counting accuracy based on the detections gathered from RetinaNet on ResNet152; ResNet50 and VGG16 backbones are 94.75%, 91.96%, and 90.96%, respectively. The speed of the counting model (Table. 5) is also reasonable compared to the number frames; e.g., video4 includes 2227 images, and the counting model only takes 5.75 seconds to run. It must be noted that this time is the processing time of just the counting model and does not include object detector processing latency.

Table 8: Counting results for all the 6 test videos. The detections are taken from the trained networks in the first fold. Extra counted means the sum of miscounted open-mouth and closed-mouth pistachios.

| Backbone Network | Ground-Truth Open-Mouth Pistachios | Ground-Truth Closed-Mouth Pistachios | Correctly Counted Open-Mouth Pistachios | Correctly Counted Closed-Mouth Pistachios | Extra Counted | Accuracy |
|---|---|---|---|---|---|---|
| ResNet152 | 409 | 152 | 397 | 145 | 11 | **0.9475** |
| ResNet50 | 409 | 152 | 386 | 152 | 24 | 0.9196 |
| VGG16 | 409 | 152 | 395 | 149 | 37 | 0.9096 |



Figure 9: Some of the examples in our dataset which are hard to classify

## 4   Discussion

Based on the Table. 7, in fold 1 and 4, ResNet152 performs better, but in other folds, ResNet50 achieves better results. On average between five-folds ResNet50 and ResNet152 get 91.87% and 91.69% mAP respectively.

It is clear that RetinaNet on ResNet models achieves better results, but the performance of ResNet50 and ResNet152 are almost equal. Although; ResNet152 is a deeper model, but in this dataset, as pistachios' features are apparent from the background and there are only two classes, there is no need to use bigger models, and the ResNet50 model is sufficient in this case.

In Fig. 9, it can be visualized that the open-mouth and the closed-mouth pistachios could look like each other, and it would be hard to distinguish them even with human eyes. So based on the detection mAP and accuracy, the models are robust for detecting the pistachios.

Our counting model's accuracy on ResNet152 backbone in fold1 is 94.75% while the detection accuracy on ResNet152 backbone in Fold1 is 86.43%. Also from Table. 8, we can see that the counting results are promising and are even higher than the detection results. This shows our proposed counter model's robustness and its ability to compensate the object detector's failures and that it can assign and track the pistachios within the video frames with high accuracy. Our counting model has been evaluated on six different videos containing 9486 frames with 561 moving pistachios and more than 350,000 single pistachios (sum of pistachios in all frames). Based on Table. 5, the counting procedure has been performed fast. This algorithm can be utilized in factories and industries related to pistachios because of its high speed and good accuracy.

## 5   Conclusion

Pistachio is a nutritious nut that originated from central Asia and the middle east and is mainly sorted into open-mouth and closed-mouth kinds. Open-mouth pistachios worth much more than closed-mouth pistachios, so it can be very valuable for the related companies to know about the exact amount of these two kinds in their packages. This paper has proposed a novel model that can be used to design a remote AI system to detect and count the number of open-mouth and closed-mouth pistachios in the factories' production line. We have introduced a new dataset that is called Pesteh-Set. It is made of 6 videos (9486 frames) and 3927 labeled open-mouth and closed-mouth pistachios.

Pistachios are motile objects and usually spin on the transportation line; therefore, from the view of a camera, open-mouth pistachios may place on their backside and appear like closed-mouth pistachios and again appear as open-mouth in some other frames. Due to this challenge, we had to develop our methods in a way to prevent false counting.

Our proposed model is constructed of a detection part and a counting part. For the detection part, we have implemented the RetinaNet object detector on our dataset to detect the different types of pistachios in video frames. We also trained and investigated RetinaNet on three different backbones: ResNet50, ResNet152, and VGG16 to find the best-performing model. The Mean Average Precision (mAP) between five-folds for RetinaNet on the ResNet50 network was 91.87%.

Our developed counter model; first receives the detected pistachios from the object detector model and then applies our designed tracking method to assign pistachios in consecutive frames. Our tracker's most crucial point is its ability to

track open-mouth pistachios that appear as closed-mouth in some frames, even in pistachios' high density and occlusion. Our counter model performs fast, with no need for GPU (besides the object detection part), and achieves good results. It was tested on six different videos containing 9486 frames with 561 moving pistachios and more than 350,000 single pistachios (sum of pistachios in all frames). This algorithm obtained 94.75% accuracy.

This work can be extended to detect and count more than two classes of pistachios, e.g., .semi-open mouth pistachios can also be added to the classes. The detection accuracy can also be increased in future works. One way is that researchers can use our developed program and the dataset to generate more labeled images. They can label all the frames of some videos and use this method [27] that was proposed to improve the motile-objects detection. Pistachios are motile objects, so by using this method, the detection accuracy should be improved.

## 6 Data Availability

In this GitHub profile (`https://github.com/mr7495/Pesteh-Set`), we have shared our dataset and all the codes that were used for preparing and labeling the dataset.

## 7 Code Availability

In this GitHub profile (`https://github.com/mr7495/Pistachio-Counting`), we made the trained neural networks, the counting algorithm and all the codes that were used for training and validating the networks, public for researchers use.

## Acknowledgment

## Conflict of Interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## References

[1] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.

[2] D. S. Berman, A. L. Buczak, J. S. Chavis, and C. L. Corbett. A survey of deep learning methods for cyber security. *Information*, 10(4):122, 2019.

[3] C. D. Burlock, G. E. Lemmons, and D. W. Williams. Apparatus for splitting closed shell pistachio nuts, Mar. 5 1991. US Patent 4,996,917.

[4] F. Chollet and Others. keras, 2015.

[5] O. Cohen, R. Linker, and A. Naor. Estimation of the number of apples in color images recorded in orchards. In *International Conference on Computer and Computing Technologies in Agriculture*, pages 630–642. Springer, 2010.

[6] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.

[7] 2018 (accessed May 2 2020).

[8] J. Feng, G. Liu, S. Wang, L. Zeng, and W. Ren. A novel 3d laser vision system for robotic apple harvesting. In *2012 Dallas, Texas, July 29-August 1, 2012*, page 1. American Society of Agricultural and Biological Engineers, 2012.

[9] Food, D. Administration, et al. Qualified health claims: Letter of enforcement discretion–nuts and coronary heart disease. docket no. 02p-0505. *Food and Drug Administration, Washington, DC*, 2003.

[10] A. Garcia-Garcia, S. Orts-Escolano, S. Oprea, V. Villena-Martinez, and J. Garcia-Rodriguez. A review on deep learning techniques applied to semantic segmentation. *arXiv preprint arXiv:1704.06857*, 2017.

[11] A. Gongal, S. Amatya, M. Karkee, Q. Zhang, and K. Lewis. Sensors and systems for fruit detection and localization: A review. *Computers and Electronics in Agriculture*, 116:8–19, 2015.

[12] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[13] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017.

[14] M. Kashaninejad, A. Mortazavi, A. Safekordi, and L. Tabil. Some physical properties of pistachio (pistacia vera l.) nut and its kernel. *Journal of Food Engineering*, 72(1):30–38, 2006.

[15] P. J. Kiger. Why pistachios are sold in their shells - unlike most nuts, Mar 2017 (accessed May 2 2020).

[16] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.

[17] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.

[18] J.-G. Lee, S. Jun, Y.-W. Cho, H. Lee, G. B. Kim, J. B. Seo, and N. Kim. Deep learning in medical imaging: general overview. *Korean journal of radiology*, 18(4):570–584, 2017.

[19] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017.

[20] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.

[21] R. Linker, O. Cohen, and A. Naor. Determination of the number of green apples in rgb images recorded in orchards. *Computers and Electronics in Agriculture*, 81:45–57, 2012.

[22] X. Liu, S. W. Chen, S. Aditya, N. Sivakumar, S. Dcunha, C. Qu, C. J. Taylor, J. Das, and V. Kumar. Robust fruit counting: Combining deep learning, tracking, and structure from motion. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1045–1052. IEEE, 2018.

[23] M. Maskan and Ş. Karataş. Fatty acid oxidation of pistachio nuts stored under various atmospheric conditions and different temperatures. *Journal of the Science of Food and Agriculture*, 77(3):334–340, 1998.

[24] H. Mureşan and M. Oltean. Fruit recognition from images using deep learning. *Acta Universitatis Sapientiae, Informatica*, 10(1):26–42, 2018.

[25] OpenCV. Open source computer vision library, 2015.

[26] M. Rahimzadeh and A. Attar. A modified deep convolutional neural network for detecting covid-19 and pneumonia from chest x-ray images based on the concatenation of xception and resnet50v2. *Informatics in Medicine Unlocked*, page 100360, 2020.

[27] M. Rahimzadeh, A. Attar, et al. Sperm detection and tracking in phase-contrast microscopy image sequences using deep learning and modified csr-dcf. *arXiv preprint arXiv:2002.04034*, 2020.

[28] M. Rahimzadeh, A. Attar, and S. M. Sakhaei. A fully automated deep learning-based network for detecting covid-19 from a new and large lung ct scan dataset. *medRxiv*, 2020.

[29] M. Rahnemoonfar and C. Sheppard. Deep count: fruit counting based on deep simulated learning. *Sensors*, 17(4):905, 2017.

[30] Different types of iranian pistachios and products of pistachios, Jan 2020 (accessed May 2 2020).

[31] O. Safren, V. Alchanatis, V. Ostrovsky, and O. Levi. Detection of green apples in hyperspectral images of apple-tree foliage using machine vision. *Transactions of the ASABE*, 50(6):2303–2313, 2007.

[32] A. Shrestha and A. Mahmood. Review of deep learning algorithms and architectures. *IEEE Access*, 7:53040–53065, 2019.

[33] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.

[34] D. C. Slaughter and R. C. Harrell. Discriminating fruit for robotic harvest using color in natural outdoor scenes. *Transactions of the ASAE*, 32(2):757–0763, 1989.

[35] D. Stajnko, M. Lakota, and M. Hočevar. Estimation of number and diameter of apple fruits in an orchard during the growing season by thermal imaging. *Computers and Electronics in Agriculture*, 42(1):31–42, 2004.

[36] J. P. Wachs, H. I. Stern, T. Burks, and V. Alchanatis. Low and high-level visual feature-based apple detection from multi-modal images. *Precision Agriculture*, 11(6):717–735, 2010.

[37] D. Whittaker, G. Miles, O. Mitchell, and L. Gaultney. Fruit location in a partially occluded image. *Transactions of the ASAE*, 30(3):591–0596, 1987.

[38] Pistachio, Apr 2020 (accessed May 2 2020).

[39] J. G. Woodruff et al. *Tree nuts: production, processing, products.* Number Ed. 2. AVI Publishing Co. Inc., 1979.