# Detection of strawberries with varying maturity levels for robotic harvesting using YOLOv4

*Zixuan He, Manoj Karkee, Priyanka Upadhyaya*

*Center for Precision and Automated Agricultural Systems, Department of Biological Systems Engineering,*

*Washington State University, Prosser, WA 99350, USA*

**Written for presentation at the**
**2021 Annual International Meeting**
**ASABE Virtual and On Demand**
**July 12–16, 2021**

**ABSTRACT.** *Inaccurate detection and/or localization of strawberries will cause fruit injury or failed attempt during robotic picking. In this work, a method is proposed to accurately detect and localize strawberries combining an object detection network, YOLOv4, and a classification network, Alexnet. YOLOv4 model was used to detect strawberries into various maturity groups (flower, immature, nearly mature, mature, and overripen berries) and provide their location information whereas Alexnet model was used for assessing if the detected matured strawberries were completely or partially visible. Compared to the same achieved by YOLOv2, YOLOv3 and YOLOv4 trained for regular objects, YOLOv4 model was specifically trained to detect small objects, which achieved the highest mean average precision of 80.7% and F1 score of 0.80 with an average precision (AP) of 91.7% in detecting mature strawberries. YOLOv4 achieved a high processing speed of 55 ms on single image (resolution: 1200×1000 pixels). The model was further validated with the RGB images generated from point cloud, where it achieved an AP of 90.15% showing that the model was robust to detect berries in images collected with different settings. Similarly, Alexnet model achieved an accuracy of 90.0% in classifying matured strawberries into completely and partially visible groups with processing speed of 3 ms per image (resolution: 227×227×3 pixels). This technique also provided 2D to 3D mapping of strawberries for the harvesting robots. This method*

*showed a strong potential as a means for providing accurate strawberry detection desirable for robotic harvesters, particularly with a collaborating dual-manipulator system.*

**Keywords.** *strawberry detection; deep neural network; YOLOv4; maturity classification; occlusion detection*

# Introduction

With delicious taste and rich nutritional characteristic, strawberry is one of most popular fruits in the United States (Hussen, A.M., 1979). Traditionally, strawberries are harvested manually, which is becoming increasingly more challenging due to increasing labor cost and continual labor shortage. To address this challenge, automated or robotic harvesting methods are being investigated around the world, which have shown to be promising alternatives to manual harvesting. One of the most important tasks of automated/robotic harvesting is to detect the fruit and their locations, and to detect maturity of agricultural product under the natural lighting conditions often characterized by uncertain and variable light intensities over both time and space.

In the past, studies used traditional methods, such as segmentation, edge detection, and feature extraction in different color spaces, to detect fruits for mechanical harvesting. Researchers have investigated the use of object shape, size (Franz et al., 1991; Woebbecke et al., 1995), and color (Yin et al., 2009) to detect fruits in specific field environments. However, it is difficult for these conventional methods to achieve desired level of accuracy and robustness in field environment with variable and uncertain lighting conditions and complex and varying cropping systems and crop canopies. Soft computing models such as support vector machine (SVM) have been used widely to address those challenges in identifying fruit (Mahendra et al., 2018). Kurtulmus et al. (2013) established a machine vision system based on SVM to detect immature peach in the branch with detecting rate of 77.9% under natural light conditions. A machine vision system was developed by Qiang et al. (2014) using SVM for tree branch and fruits identification in field conditions and achieved an accuracy of 92.4%. The model was trained using RGB images, which achieved better results than the same achieved by threshold-based methods but was easily affected by outdoor lighting conditions. Contrast limited adaptive histogram equalization (CLAHE), developed by Choi et al. (2016), was used for improving the illumination of image features to detect dropped citrus fruit on

the ground and to evaluate the decay stages, which achieved an accuracy of 89.5%. However, outdoor lighting conditions still presented challenges on the fruit detection task. These traditional machine learning methods made improvement in detecting agricultural objects like fruit but were still found to lack flexibility and reliability when lighting and other environmental conditions varied (Kamilaris & Prenafeta-Boldú, 2018).

With the introduction of deep learning, one of the latest machine learning techniques developed, to computer vision tasks in agriculture in recent years, the performance of machine vision systems (e.g. accuracy and robustness) were improved. With the rapid advancement in computational hardware including powerful graphical processing units (GPUs), deep convolutional neural networks (DCNN) have been widely used for detection and grading on agricultural products. CNN-based systems were increasingly used in developing robotic harvesting machines, which helped robots to accurately locate fruits and/or other types of corps and estimate fruit/crop maturity.  Lamb et al. (2018) developed a strawberry detection system using convolutional neural networks (CNN) and reported an average precision (AP) of 0.877 at 360×640 pixels resolution with a computational speed of 0.23 frames-per-second under natural lighting conditions. Region-based convolutional neural networks (R-CNN), whose mean average precision was 53.7% on PASCAL 2010 dataset (Everingham et al., 2010), combine region proposals, such as region boxes, with CNNs (Gishick et al., 2014).  Using this network design, Christiansen et al. (2016) developed a real-time machine vision system (called DeepAnomaly) to detect the obstacles, like people, animals, and other obstacles in agricultural fields. Chen et al. (2019) applied a Faster R-CNN with backbone of ResNet50 (He et al., 2016) to conduct yield prediction by counting the number of strawberries in the field. Compared with R-CNN, the model achieved a reduced training time and improved processing speed and reported a mean accuracy of 84.1% on average. However, it is still difficult for these methods based on two stage networks including dense prediction and sparse perdition to meet the requirement of real-time detection although they achieved a relatively high accuracy and robustness in detecting fruit in variable lighting conditions.

Compared to the two-stage networks like R-CNN, Faster-RCNN, and Mask-RCNN, You-Only-Look-Once (YOLO) models (Redmon et al., 2016; Redmon et al., 2017; Redmon, et al., 2018;  Bochkovskiy, et al., 2020) put the detection results including bounding boxes and class probability directly with a single feed forward network, which makes YOLO models computationally much more efficient. YOLOv2 was improved greatly on detection accuracy and learning process from YOLO when anchor was used in YOLOv2, which was inspired from faster-RCNN (Redmon et al., 2017). YOLOv3 was built on YOLO and YOLOv2 but included a specific feature (binary cross-entropy loss) for each label in training dataset (Redmon, et al., 2018). YOLOv3 become preferred network for object detection in various applications as its computational performance was superior allowing machine vision systems to achieve real time object (e.g. fruit) detection. Researchers have achieved promising results in fruit detection adopting YOLOv3 structure (Liu, et al., 2020; Lawal et al., 2021; Yu, et al., 2020). However, these studies have mostly focused on detecting only two or three classes of big objects which were over

40×40 pixels in images with a resolution of 500×500 pixels or more. In addition, with introduction of Spatial pyramid pooling (He et al., 2015) and Path aggregation network (Liu et al., 2018), it was demonstrated that YOLOv4 achieved an AP of 43.5% for the MS COCO dataset (Lin et al., 2014) and ~65 FPS processing time on Tesla V100 GPU of YOLOv4 (Bochkovsiy, et al., 2020), which was an improvement of 10% and 12%, respectively, compared to the same with YOLOv3. However, there are limited studies in applying YOLOv4 for fruit detection in field conditions, especially for small fruit/object detection.

Strawberries are easily damaged and bruised during mechanical/robotic harvesting, and therefore machine vision systems used in these machines need to be highly accurate in detecting and localizing berries. Yang et al. (2019) used Mask-RCNN to detect mature and immature groups of strawberries with a high recognition rate of 98.4% but the processing speed was only 8 frame per second, on average. In addition, two class classification has a potential to cause more errors for berries in the borderline state between mature and immature ones, and there is no provision to find overripen berries. Lab fruit, a fruit detection system developed by Kirk et al. (2019), applied one-stage network to achieve real-time strawberry detection from different shooting angles to reduce the error caused by occlusion. Yang et al (2020) developed a machine vision system based on modified YOLOv3 for identifying only mature strawberries and vines and achieved a success rate of 84.4%. Although these past studies achieved good performance on strawberry detection, they focused on only two or three levels of strawberry maturities, which can lead to less precise harvest decision (which fruit is optimal for harvest). These past studies didn't have high accuracy on localization or computational speed for real-time application in robotic harvesting systems.

The specific objectives of this study were:

1) To develop a machine vision system for detecting and localizing flowers and strawberries of four different maturity levels (flower, immature, nearly mature, mature, overripen) in field conditions based on a YOLOv4 model.

2) To develop a technique to classify mature (optimal targets for harvesting) berries (detected in Obj# 1) to completely visible and partially visible/occluded groups.

The results from this multi-class detection and classification technique will be highly valuable to develop an efficient harvesting robot, particularly with dual- or multi-manipulator system that can avoid interfering with the canopy areas where immature and nearly mature berries are present, and avoid going to overripen berries for picking. The information can also be useful to allow collaborative manipulators to work together to open up the canopies in the area where partially visible, mature berries are detected so that they could be better exposed for precise picking while directly going for picking the completely visible berries thus improving the efficiency and productivity.

# Material and Methods

**Data Collection**

Images used in this experiment were collected in a commercial strawberry field (located at 28°41'47.4"N，81°12'14.4"W; Orlando, Florida), shown in figure 1. The normal strawberry growing season in Florida is usually between December and the following April.  ZED2 stereo camera was used for acquiring strawberry canopy images (specification in table 1). The camera was mounted on a mobile platform for capturing canopy images including RGB images, depth images and point cloud files. Data collection lasted from February 1, 2020 to February 22, 2020; the images ware collected twice a day, around 10:00 a.m. – 12:00 a.m. and 2:00 p.m. – 5:00 p.m.  The canopy images were collected from ~100 cm above the plant bed using the nadir view of the camera. More than 5,000 RGB images, and 1,000 depth images and point cloud files of strawberry canopies were collected (example data in figure 2).



**Figure 1. A commercial strawberry field at 28°41'47.4"N，81°12'14.4", Orlando, Florida**

**Table 1. General Specifications of ZED2 cameras**

| Item | Parameters |
|---|---|
| Output Resolution | Side by Side 2x (2208x1242) @15fps |
|  | 2x (1920x1080) @30fps |
|  | 2x (1280x720) @60fps |
|  | 2x (672x376) @100fps |
| Field of View | Max. 110°(H) x 70°(V) x 120°(D) |
| Interface | USB 3.0/2.0 - Integrated 1.2m cable |
| Depth Range | 0.3 m to 20 m (1 to 65.6 ft) |
| Depth Accuracy | < 1% up to 3m |
|  | < 5% up to 15m |

**(a)**



**(b)**



**(c)**

**Figure 2. Example of the data collected and used in this study: (a) an RGB image, (b) a depth image, and (c) a point cloud data.**

## Data Pre-processing

Strawberries in the vine canopies undergo a series of physiological and biochemical changes during the growing season, which leads to the changes in fruit color and other characteristics unique to different maturity levels of berries (Caiqin et al., 1998). In this study, maturity development of strawberries was divided in five categories (flower and four maturity stages of fruits). The fruit maturity levels were defined based on the fruit grading method described by Barnes et al. (1976) and as shown below:

  Immature stage – small and large green berries

 Nearly mature stage – reddish berries (1/ 3 to 3/4 of the fruit surface area covered by red color)

  Mature stage – red ripe berries (red color over 3/4 of the surface area)

  Overripen stage – senescent berries

To detect objects (flowers and strawberries with four different maturity levels) using YOLOv4 model, the location, and maturity of individual strawberries in each canopy image need to be labelled. The RGB images on strawberry canopies were cropped to remove unwanted part of the background and resized to 1200×1000 pixels. Annotation/labelling on strawberries in five individual classes was performed with rectangle bounding boxes (see figure 3) using an open-source graphical image annotation tool, labelimg_v1.8.0 (Amazon, Seattle, WA). Bounding box data contained a five-element vector of the form [class x/1000 y/1200 width/1000 height/1200]. This vector specifies the class label for the object, the center point and size

of the corresponding bounding boxes.



**Figure 3. An example image with bounding box labels generated using labelimg_v1.8.0 (Amazon, Seattle, WA)**

## YOLOv4-based object detection

*YOLOv4 model structure*

The 4th generation of YOLO (YOLO v4) was developed in May 2020. YOLOv4 consists of backbone (CSPDarknet 53), neck (Spatial pyramid pooling and Path aggregation network), and a head (YOLOv3). The structure of YOLOv4 is shown in figure 4.



Input (Strawberry canopy images)

Backbone(CSPDarknet 53)

Neck(SSP,PAN)

Dense prediction (YOLO v3)

Strawberry detection

**Figure 4. A generic architecture of YOLOv4 model including backbone (CSPDarknet53), neck (SSP and PAN), and dense prediction (YOLOv4)**

CSPDarknet 53 convolution network, as the backbone of YOLO v4, is responsible for feature extraction, which used a CSPNet strategy (Wang et al., 2020) to partition the feature map of the base layer into two parts and then merge them through a cross-stage hierarchy while keeping most of the structure of Darknet 53 (Redom & Farhadi, 2018). The use of a split and

merge strategy allowed for more gradient flow through the network. The whole structure of CSPDarknet 53 is shown in figure 5(a). CSPDarknet 53 has 53 layers, which includes 3*3 convolution layers and 1*1 convolution layers. The convolution layer is composed of a batch normalization layer and LeakyReLu layer. Following the convolution layers, a residual layer is appended, which is also composed of 2 convolution layers, as shown in figure 5(b). The residual layers were added to the structure of CSPDarknet 53 to address the problem of the gradient disappearance or gradient explosion. CSPDarknet53 (Wang et al., 2020) was found to be better when used in detecting objects on the MS COCO dataset (Lin, T. Y. et al., 2014) compared to the same with CSPResNet50 (Wang et al., 2020). CSPDarknet 53, as a backbone of YOLO v4, has ability to detect five classes of strawberries in this study.
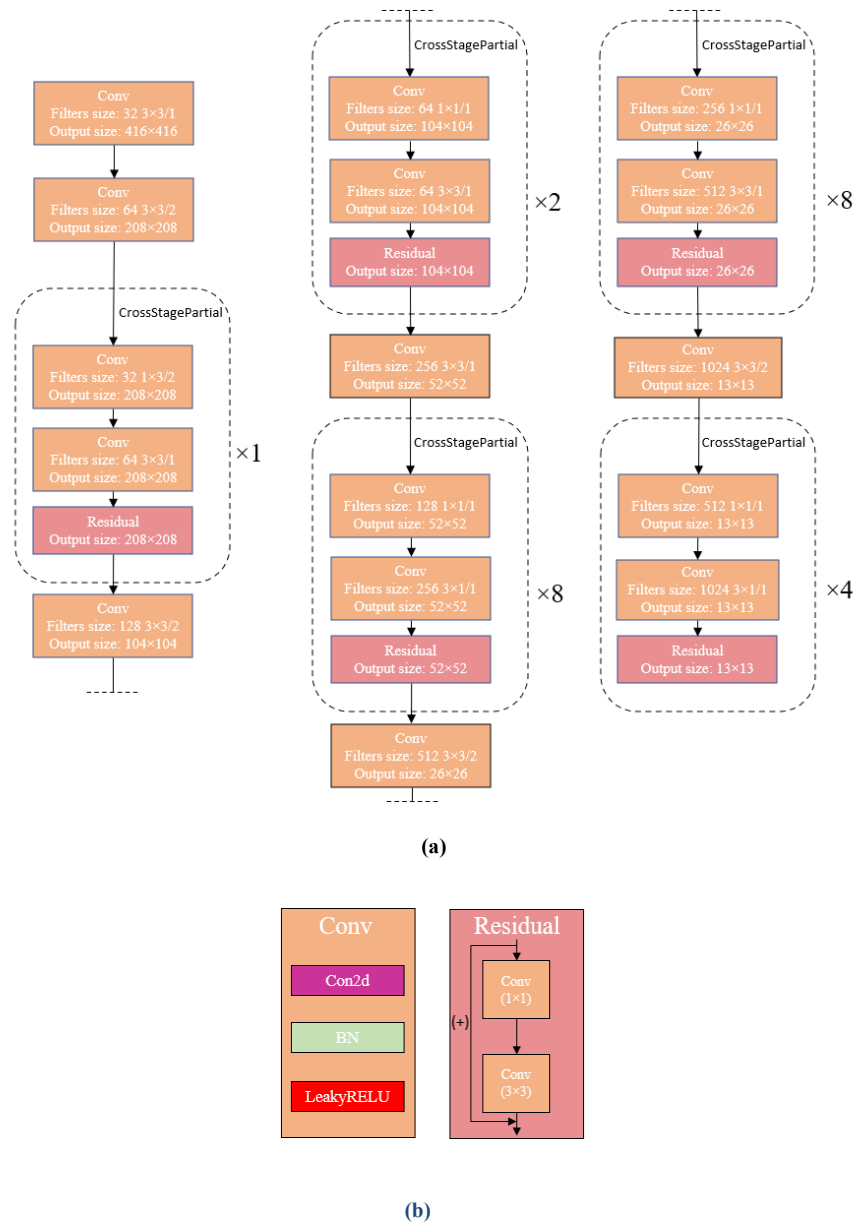


**(a)**



**(b)**

**Figure 5. A typical network architecture of CSPDarknet 53: (a) CSPDarknet 53 structure with 53 Layers, and (b) structure of convolution layers and residual layers built on 2 convolution layers.**

The neck of YOLOv4 is composed of additional blocks called SPP and PAN. The SPP was added over CPSDarknet 53,

which greatly separated out the most important feature of context and didn't put any load on processing speed. Some previously modified YOLO models (Lawal et al., 2021) have used this method over YOLOv3 model to achieve increasing receptive field effectively. PAN was used as a method of parameter aggregation from different backbone levels for different detector levels to replace feature pyramid networks (FPN) used in YOLOv3.

*Model training for small object detection*

YOLO v4 model was structured such that it received RGB images with resolution of 640×768 pixels as inputs. The learning rate, which decreased gradually during the training progress, was set to be 0.001 in the first 8000 iterations, 0.0001 between 8001 and 9000 iterations, and 0.00001 between 9001 and 10000 iterations. Decreasing learning rate was used to help reduce training loss. The batch subdivision of 64 was used to reduce the memory usage during training. Most of the labelled strawberries were small occupying <50×50 pixels in RGB image with pixel resolutions of 1200×1000 pixels. Because the target objects to be detected were small, it was necessary to set the YOLOv4 model for detecting small objects. YOLOv4 model can have better performance on detecting small objects through improving the ability of extracting the feature map of small objects when below settings were applied.

(a) The number of layers of the sixth route layers after SPP was changed to 23 from 54 in regular setting.

(b) The stride in upsample layer before sixth route layer after SPP was changed to 4 from 2 in regular setting.

(c) The stride in convolutional layer after first yolo layer was changed to 4 from 2 in regular setting.

The momentum and decay for the network were set to be 0.949 and 0.05 respectively. A dataset including 1,200 regular RGB images and 200 RGB images generated form point cloud data was used to train the network. A test dataset, composed of 100 images of strawberry canopies, was used to assess the performance of YOLO v4 model during training. The model was trained in a laptop computer with NVIDIA GeForce GTX 1070 GPU and intel core i7-8750H CPU (2.20 GHz speed and 32.0 GB RAM).

To assess the performance of this model with the same achieved by previous versions, YOLOv2, YOLOv3 and YOLOv4 (normal setting) were also implemented and trained with same dataset. Training and test results with these models were then compared against the results achieved by YOLOv4 (small object) model.

**Classifying mature strawberries to completely and partially visible groups**

AlexNet, a deep convolutional neural network consisting of 8-layer CNN including 5 convolution layers and 3 fully connected layers (Iandola, F. N., et al, 2016), was used to classify mature strawberries (those within the bounding boxes) detected by YOLOv4 into completely visible and partially visible groups. Bounding boxes of completely visible strawberries would be saved while bounding boxes of partially visible strawberries would be ignored during providing location of mature

strawberries after classification of Alecnet.

The size of input images to AlexNet was 227×227×3 pixels. After the detection by YOLOv4 model, a dataset including bounding boxes, labels and probabilities were created by cropping matured strawberries out of the images and resizing them to meeting the input image resolution requirement (227x227 pixels; Figure 6). The images were labeled manually to assign them into completely (430 images) or partially visible (367 images) classes. The AlexNet model was trained with those images in the same laptop computer used for YOLOv4 model training. The model was trained for 20 epochs with a learning rate of 0.0001. Training dataset included 70% of the images and test dataset included the rest of images.



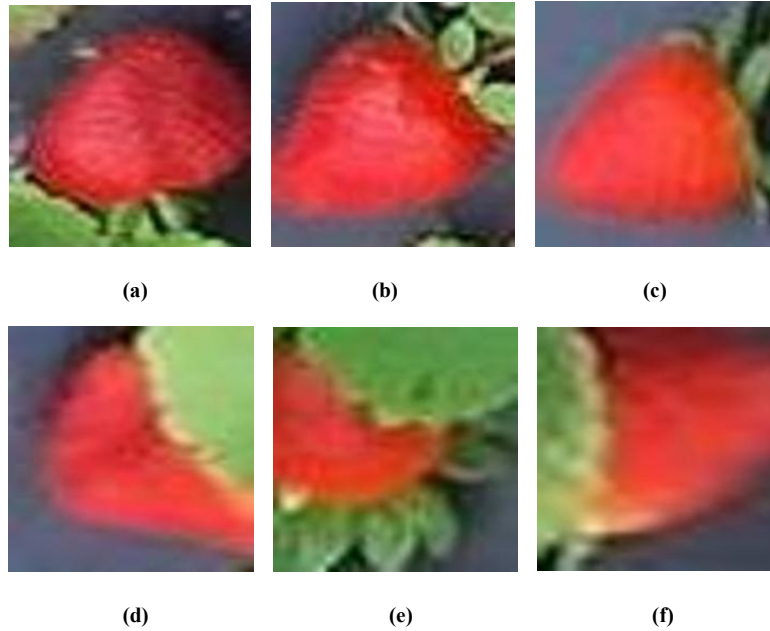|  (a)  |  (b)  |  (c)  |
|  (d)  |  (e)  |  (f)  |

**Figure 6. Selected examples of mature strawberry images cropped from bounding boxes: (a)(b)(c) complete strawberries under bounding boxes, (d)(e)(f) overlapped/partially visible strawberries under bounding boxes.**

## Performane Assessment

Strawberry detection results were evaluated using recall (R), precision (P), F1-score and average precision (AP) and mean value of average precisions for each class (mAP) with intersection-over-union (IOU) of 50%.

The IOU is used for evaluating the correctness of bounding box on detected objects. The IOU is calculated as follows:

$$IOU = \frac{|A \cup B|}{|A \cap B|} = \frac{Area(I)}{Area(U)} \tag{1}$$

Where A is predicted bounding boxes, B is ground truth bounding boxes, Area(I) is the intersection of predicted and ground truth bounding boxes, and Area(U) is the union area of predicted and ground truth bounding boxes (Fu et al., 2018).

P, R, and F1 score are used for measuring the accuracy of overlap between predicted and ground truth bounding boxes, which are calculated as follows:

$$Precision = \frac{TP}{TP+FP} \tag{2}$$

$$Recall = \frac{TP}{TP+FN} \tag{3}$$

$$F1\ score = 2 \times \frac{Precison \times Recall}{Precision + Recall} \tag{4}$$

Where TP is the number of true positive strawberries detected, FP is the number of false strawberries detected, and FN is the number of strawberries falsely not detected as strawberries. The IOU threshold was set at 0.5, which means that the result is true positive if the IOU between the predicted result and ground truth bounding boxes was greater than 0.5. F1 score was used for showing the comprehensive performance of the models. Average Precision (AP) was used for evaluating performance of the model in different classes and mAP was adopted to show the overall performance under different confidence thresholds. AP and mAP can be calculated as follow:

$$AP = \sum_n (r_{n+1} - r_n) \max_{\tilde{r}:\ \tilde{r}^3 r_{n+1}} p(\tilde{r}) \tag{5}$$

$$mAP = \frac{1}{N}\sum_{i=1}^{N} AP_i \tag{6}$$

Where $p(\tilde{r})$ is the Precision at Recall $\tilde{r}$.

Classification of completely and partially visible mature strawberries was evaluated using accuracy, which can be calculated as follows:

$$Accuracy = \frac{number\ of\ True\ Objects\ in\ the\ Class}{number\ of\ Total\ Objects\ in the\ Class}$$

$$\tag{9}$$

# Results

**Multi-class Strawberry Detection**

Performance of the trained models was evaluated qualitatively and quantitatively using a canopy image dataset collected in different time periods in a commercial field. Both regular RGB images and RGB images generated from point cloud (P-RGB images) were used to test the performances of the trained YOLOv4 model.

*Performance of YOLOv4 in detecting strawberries with different maturity levels*

The performance in detecting strawberries in each maturity group was assessed with AP (Henderson, 2010) and F1 score. The performance measures including AP for each class, mAP and IOU achieved by the YOLOv4 (trained for small object) model with varying iterations of training are shown in Table 2. The performance was evaluated using the test dataset consisting of 100 images of strawberry canopies acquired in the field conditions.

**Table 2. Performance of YOLOv4 (trained to detect small objects) model in detecting strawberry flowers and berries at different maturity levels**

| Training Iterations | AP of Individual Classes (%) | | | | | mAP(%) | IOU(%) |
|---|---|---|---|---|---|---|---|
| | Flower | Immature | Nearly mature | Mature | Overripen | | |
| 1000 | 8.60 | 62.84 | 20.49 | 76.87 | 18.04 | 39.05 | 37.84 |

| | | | | | | |
|---|---|---|---|---|---|---|
| 2000 | 58.70 | 84.04 | 73.17 | 87.99 | 51.03 | 70.99 | 54.00 |
| 3000 | 58.62 | 83.58 | 75.63 | 87.25 | 42.85 | 69.59 | 52.45 |
| 4000 | 59.35 | 85.49 | 79.45 | 87.25 | 44.83 | 71.27 | 51.52 |
| 5000 | 66.30 | 89.12 | 80.33 | 90.03 | 58.29 | 76.81 | 57.40 |
| 6000 | 66.79 | 88.15 | 84.17 | 90.92 | 66.24 | 79.26 | 58.45 |
| 7000 | 67.17 | 87.71 | 85.28 | 91.73 | 71.56 | 80.69 | 58.87 |
| 8000 | 70.66 | 88.25 | 84.25 | 91.12 | 68.99 | 80.71 | 57.01 |
| 9000 | 68.31 | 86.14 | 83.34 | 88.85 | 62.62 | 77.85 | 53.98 |
| 10000 | 70.82 | 88.62 | 84.10 | 87.98 | 61.08 | 78.52 | 47.78 |

From table 2, it was found that the YOLOv4 model was gradually improving until 8,000 iterations (mAP of 80.71%) but then slightly deteriorated (by 2.19%) in later 2,000 iterations. The results showed that the performance of YOLOv2 model in detecting immature fruits, nearly mature fruits, and mature fruits, which are the core classes of strawberry detection, is better than in detecting flowers and overripen fruit. The main reason of the different level of detection performance is the number samples of each class available in the training dataset. The number of flowers and overripen fruits in training dataset was substantially smaller compared to the number of immature, nearly mature, and mature fruit. The AP for mature fruit was achieved over 90% when the model was trained to 7,000 to 8,000 iterations. As matured berries are the objects of interest for harvesting, the good performance with this class is essential even with a slightly reduced performance with the flowers and overripen fruit. In addition, the AP for nearly mature class was over 84% when the model was trained to 7000 to 8000 iterations, which helped to reduce the error in determining if a strawberry was mature or immature, and in determining which berries are ready for harvesting. Occlusion was one of major reasons causing errors in detecting strawberries in all classes. YOLOv4 model, as with any other fruit detection models using RGB imaging, can only detect visible berries (at least some part of the object visible). In addition, for some partially visible berries, as only the visible parts can be detected, can lead to positioning error during strawberry detection. Another challenge in classification was caused by nearby parts of calyx or leaves overlapping with given fruit, which made it challenging to distinguish mature berries and nearly mature berries.

The proposed method achieved good performance in detecting nearly mature, and mature strawberries whereas the performance was not as good with the flower and overripen groups. The physical structures of flowers and overripen fruits are more complex than structures of other three classes and were somewhat similar to the leaves or dead leaves, which could have caused some issues in detecting them correctly with YOLOv4 (small object) model. Increasing the dataset size could be helpful in the future to improve AP for all the strawberry classes. Additionally, more than two cameras could be applied from different angles to detect strawberries, which can provide better point cloud of canopies and also help minimize occlusion problems in the canopy.

*Comparison of YOLO models on strawberry detection*

As mentioned before, different versions of YOLO models were trained and tested for strawberry detection and to compare their performances. Fgure 7 shows some example images with different classes of objects detected and visualized. The red and blue bounding boxes in the figure were labelled manually to represent the missed detection (FN) and false detection (FP) in different classes. From the figure 7, it can be viewed, qualitattively, that the YOLOv4 (small objects) achieved the best performace while there are 2 missed detction with YOLOv2 (figure 7a), and 2 missed detection and 1 false detection with YOLOv3 (figure 7b). Although there was no substantial differene in F1 score and mAP (table 3, discussed later) between YOLOv4 and YOLOv4 (small objects), there was one missed detction and one false detction with YOLOv4 (figure 7c).



(a)

(b)

(c)

(d)

**Figure 7. Selected example images with depiction of missed detections (red boxes) and wrong detections (blue boxes) of strawberries in different groups with (a)YOLO v2, (b)YOLOv3, (c) YOLOv4, and (d) YOLOv4 (small objects).**

The quantitative performance measures based on training dataset and test dataset were shown in table 3 and 4. The average processing time of YOLO models on single image with the resolution of 1200×1000 pixels was shown in table 5.

*Table 3. Training performance of various YOLO models on detecting strawberry classes*

| Methods | AP of Class (%) | | | | | mAP (%) | F1 |
|---|---|---|---|---|---|---|---|
| | Flower | Immature | Nearly mature | Mature | Overripen | | |
| YOLOv2 | 55.46 | 70.88 | 82.73 | 84.87 | 67.08 | 72.20 | 0.73 |
| YOLOv3 | 79.25 | 91.80 | 86.48 | 93.74 | 76.94 | 85.64 | 0.82 |
| YOLOv4 | 86.84 | 94.45 | 95.95 | 96.73 | 89.31 | 92.66 | 0.86 |

| | | | | | | |
|---|---|---|---|---|---|---|
| YOLOv4(small objects) | 88.71 | 97.23 | 93.21 | 97.23 | 92.29 | 93.24 | 0.86 |

**Table 4. Performance of various YOLO models in detecting strawberry classes based on the test dataset**

| Methods | AP of Class (%) | | | | | mAP (%) | F1 |
|---|---|---|---|---|---|---|---|
| | Flower | Immature | Nearly mature | Mature | Overripen | | |
| YOLOv2 | 45.66 | 65.27 | 70.09 | 79.04 | 44.60 | 60.93 | 0.65 |
| YOLOv3 | 65.00 | 84.00 | 65.67 | 85.44 | 59.93 | 72.01 | 0.73 |
| YOLOv4 | 67.75 | 86.01 | 83.82 | 88.82 | 67.95 | 78.87 | 0.78 |
| YOLOv4 (small objects) | 71.51 | 87.71 | 85.28 | 91.73 | 68.99 | 80.69 | 0.80 |

**Table 5. Average processing time of YOLO models for an image with resolution of 1200×1000 pixels**

| Model | Time (ms) |
|---|---|
| YOLOv2 | 30.80 |
| YOLOv3 | 50.75 |
| YOLOv4 | 55.26 |
| YOLOv4(small object) | 55.19 |

The results (tables 3 and 4) showed that YOLOv2 didn't perform well on both training and test datasets achieving the lowest mAP of 60.9% with test dataset among the YOLO versions evaluated. Meanwhile, the performance of YOLOv3 increased substanially compared to YOLOv2, as shown by 11.08% increase in AP with the test dataset. When it comes to YOLOv4 and YOLOv4 (specialized to small objects), it was found that both models were substantially better in detecting various strawberry classes compared to YOLOv2 and YOLOv3, with mAP over 78.00% in test dataset and over 92.00% in training dataset. When comparing YOLOv4 (small object) with YOLOv4 (normal) in detecting mature strawberries, it was found that mAP of YOLOv4 (small objects) increased by 1.82% from YOLOv4 while F1 score was increased by 0.02 (marginal improvements). The AP achieved by YOLOv4 (small object) in each class were also slightly higher than AP achieved by YOLOv4 (table 4), in general. It was found that AP of nearly mature class achieved by YOLOv4 on training dataset was slightly higher (by 3.74%) than the same achieved by YOLOv4 (small object) model, which, however, was

higher with the YOLOv4 (small object) model when it comes to the test dataset.

From table 5, it was found that YOLOv2 was the fastest model among the tested YOLO versions in detecting strawberries classes on a given image (resolution of 1200×1000 pixels), which was expected as YOLOv2 has a much simpler architecture compared to YOLOv3 and YOLOv4. Compared to YOLOv2, the processing speed of YOLOv3 increased by 19.95 ms while YOLOv4 and YOLOv4 (small objects) have slower speeds of 55.26 ms and 55.19 ms, respectively. As can be seen and as expected, there was no noticeable difference on processing speed of YOLOv4 and YOLOv4 (small objects).

*Detection Results of mature strawberry on Point Cloud*

The point clouds, which were dense and complete, were collected from 100 cm above the plant bed using ZED2 camera. The P-RGB images from point cloud had gaps in the both sides of the plant bed, which are the unnecessary parts and therefore were removed before strawberry detection but the canopies in the center of the bed were constructed completely. Examples of P-RGB images and RGB image were shown in figure 8.
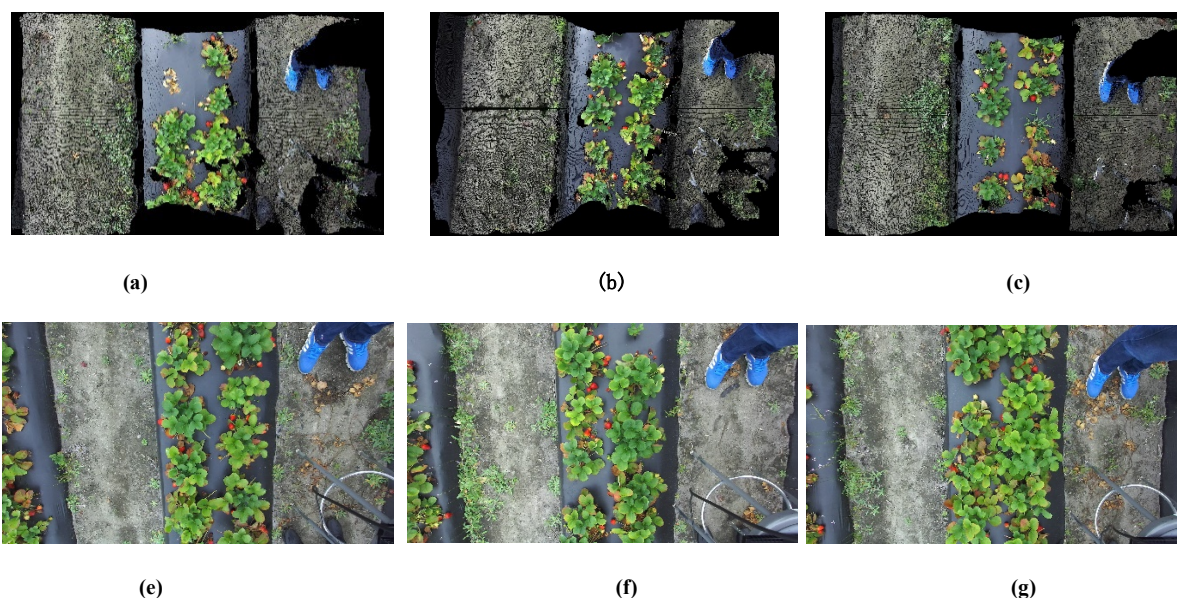


| (a) | (b) | (c) |



| (e) | (f) | (g) |

**Figure 8. Examples of P-RGB images and RGB images: (a)(b)(c) examples ofRGB images generated form point cloud, and (e)(f)(g) regular RGB images.**
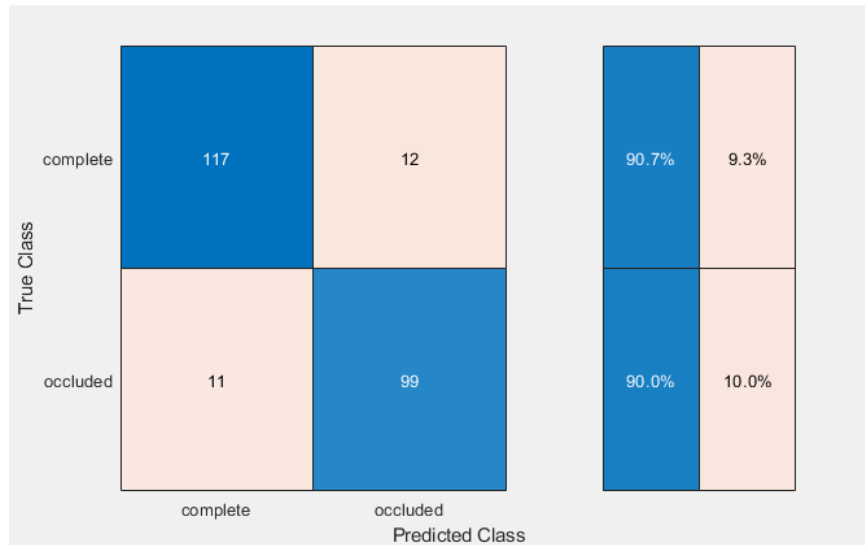
The performance of YOLOv4 (small object) was assessed using 50 P-RGB images. The number of P-RGB images used in training dataset was only a small fraction of the total RGB images used. With this test dataset, the model achieved an AP of 90.5% in detetcting mature berries, which was slightly lower than when the model was tested with normal RGB images (AP of 91.73%). This result indicated that the model is robust to process images collected with different kinds of sensors and setting. However, more comprehesive research would be necessary to fully validate this finding. The labels of matrue strawberries detected on P-RGB images were projected on the point clouds (figure 12). The location information of mature starwberries then can be acquired from these point coulds and can be sent to the controller for robotic harvesting.
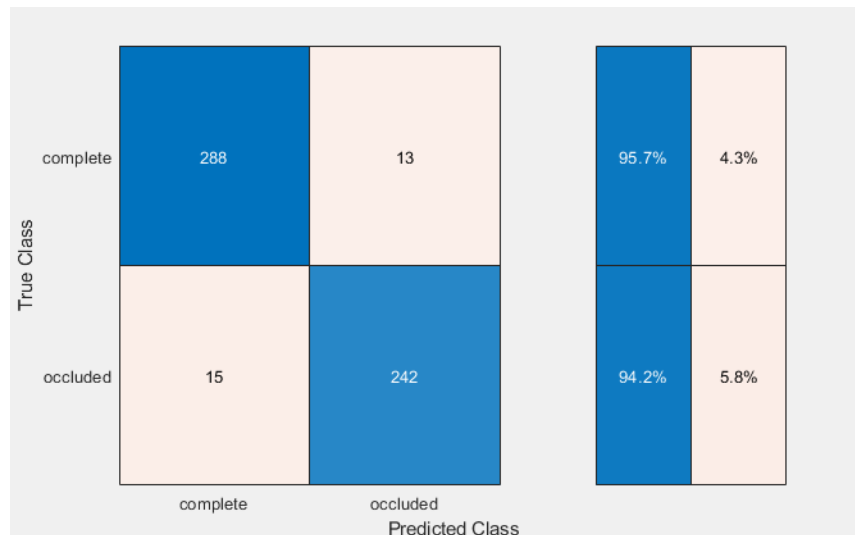
**Figure 9. Selcected examples of completely visible strawberry detection on point cloud based on YOLOv4 and Alexnet**

**Classification of mature berries to completely and partially visible categories**

From figure 10(a) and figure 10(b), it was found that the accuracy of AlexNet in detecting completely visible berries was 90.7% and the same in detecting partially visible berries was 90.0% with the test dataset. The performance of the model was slightly better with the training dataset (95.7% for completely visible and 94.2% for partially visible berries. Completely visible group had 60 more images than overlapped group, which might have caused the slight difference in the performance with two individual classes. The processing time per image (227×227×3 pixels) for this model was 3ms on average.



**(a)**

**Figure 10.   Confusion matrix of dataset from Alexnet on complex group and occluded group: (a)test dataset, (b) training dataset.**

# Conclusion

In this paper, we presented a machine vision system based on YOLOv4 for strawberry detection including flowers and those in four different maturity levels. A post-processing method base on Alexnet was also proposed to classify mature berries into completely and partially visible groups for improving the detection result on mature class. This system had ability to provide the presence, location, and maturity information of strawberries (including visibility information of mature berries) in the canopy under field condition.

A ZED2 camera was used to collect RGB images, depth images and point clouds of strawberry canopies from 100 cm above the plant beds. The performance of YOLOv4 (small object) model was evaluated using average precision (AP), mean average precision(mAP), and intersection of union (IOU). With a test dataset of 100 images, a mAP of 80.7% with IOU of 57.0% was achieved in detecting mature berries when the model was trained to 8,000 iterations. The model took 55 ms to process a single RGB image with a resolution of 1200×1000 pixels. It was found that, among the classes studied, overripen fruit were the most difficult to detect (AP of 69.0%) as the overripen fruit had multiple structures and the model may mix them with dead leaves. The performance of different YOLO models was also compared, which showed that YOLOv4 (small object) had the best performance in strawberry detection task with highest F1 score of 0.80. Besides, the accuracies of AlexNet model was found to be 90.7% in classifying mature berries to completely visible category and 90.0% in classifying to partially visible category. This information can be useful to allow collaborative manipulators to work together to open up the canopies in the area where partially visible, mature berries are found so that they could be better exposed for precise picking while directly going for picking the completely visible berries thus improving the efficiency and productivity.

As mentioned before, the model did not perform well in detecting flower and overripen berries Increasing the dataset size could be helpful in the future to improve AP for all the strawberry classes. Additionally, more than two cameras could be applied from different angles to detect strawberries, which can provide better point cloud of canopies and help minimize occlusion problems in the canopy. In addition, reconstruction of partially visible, mature strawberries and estimation of corrected berry locations could be investigated in the future.

# References

Barnes, M. F., & PATCHETT, B. J. (1976). Cell wall degrading enzymes and the softening of senescent strawberry fruit. *Journal of food science*, *41*(6), 1392–1395.

Bochkovskiy, A., Wang, C. Y., & Liao, H. Y. M. (2020). Yolov4: Optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934*.

Caiqin, G., & Dongxue, Z. (1998). The physiologic and biochemical changes of strawberry during maturation. *Journal of Mountain Agriculture and Biology*, *17*(6), 345–348.

Chen, Y., Lee, W. S., Gan, H., Peres, N., Fraisse, C., Zhang, Y., & He, Y. (2019). Strawberry yield prediction based on a deep neural network using high-resolution aerial orthoimages. *Remote Sensing*, *11*(13), 1584.

Christiansen, P., Nielsen, L. N., Steen, K. A., Jørgensen, R. N., & Karstoft, H. (2016). DeepAnomaly: Combining background subtraction and deep learning for detecting obstacles and anomalies in an agricultural field. *Sensors*, *16*(11), 1904.

Choi, D., Lee, W. S., Ehsani, R., Schueller, J., & Roka, F. M. (2016). Detection of dropped citrus fruit on the ground and evaluation of decay stages in varying illumination conditions. *Computers and Electronics in Agriculture*, *127*, 109–119.

Everingham, M., Van Gool, L., Williams, C. K., Winn, J., & Zisserman, A. (2010). The pascal visual object classes (voc) challenge. *International journal of computer vision*, *88*(2), 303-338.

Franz, E., Gebhardt, M. R., & Unklesbay, K. B. (1991). The use of local spectral properties of leaves as an aid for identifying weed seedlings in digital images. *Transactions of the ASAE*, *34*(2), 682–0687.

Girshick, R., Donahue, J., Darrell, T., & Malik, J. (2014). Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 580–587).

He, K., Zhang, X., Ren, S., & Sun, J. (2015). Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE transactions on pattern analysis and machine intelligence*, *37*(9), 1904–1916.

Hussen, A. M. (1979). Estimated costs and returns from mechanical strawberry harvest in Oregon: a progress report.

Lawal, M. O. (2021). Tomato detection based on modified YOLOv3 framework. *Scientific Reports*, *11*(1), 1–11.

Liu, G., Nouaze, J. C., Touko Mbouembe, P. L., & Kim, J. H. (2020). YOLO-tomato: A robust algorithm for tomato detection based on YOLOv3. *Sensors*, *20*(7), 2145.

Kamilaris, A., & Prenafeta-Boldú, F. X. (2018). Deep learning in agriculture: A survey. *Computers and electronics in agriculture*, *147*, 70–90.

Kurtulmus, F., Lee, W. S., & Vardar, A. (2014). Immature peach detection in colour images acquired in natural illumination conditions using statistical classifiers and neural network. *Precision agriculture*, *15*(1), 57–79.

Lamb, N., & Chuah, M. C. (2018, December). A strawberry detection system using convolutional neural networks.

In *2018 IEEE International Conference on Big Data (Big Data)* (pp. 2515–2520). IEEE.

Lawal, M. O. (2021). Tomato detection based on modified YOLOv3 framework. *Scientific Reports*, *11*(1), 1–11.

Lin, T. Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., ... & Zitnick, C. L. (2014, September). Microsoft coco: Common objects in context. In *European conference on computer vision* (pp. 740–755). Springer, Cham.

Liu, G., Nouaze, J. C., Touko Mbouembe, P. L., & Kim, J. H. (2020). YOLO-tomato: A robust algorithm for tomato detection based on YOLOv3. *Sensors*, *20*(7), 2145.

Lin, T. Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., ... & Zitnick, C. L. (2014, September). Microsoft coco: Common objects in context. In *European conference on computer vision* (pp. 740–755). Springer, Cham.

Liu, S., Qi, L., Qin, H., Shi, J., & Jia, J. (2018). Path aggregation network for instance segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 8759–8768).

Mahendra, O., Pardede, H. F., Sustika, R., & Kusumo, R. B. S. (2018, November). Comparison of features for strawberry grading classification with novel dataset. In *2018 International Conference on Computer, Control, Informatics and its Applications (IC3INA)* (pp. 7–12). IEEE.

Qiang, L., Jianrong, C., Bin, L., Lie, D., & Yajing, Z. (2014). Identification of fruit and branch in natural scenes for citrus harvesting robot using machine vision and support vector machine. *International Journal of Agricultural and Biological Engineering*, *7*(2), 115–121.

Kirk, R., Cielniak, G., & Mangan, M. (2020). L* a* b* fruits: A rapid and robust outdoor fruit detection system combining bio-inspired features with one-stage deep learning networks. *Sensors*, *20*(1), 275.

Redmon, J., Divvala, S., Girshick, R., & Farhadi, A. (2016). You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 779–788).

Redmon, J., & Farhadi, A. (2018). Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*.

Redmon, J., & Farhadi, A. (2017). YOLO9000: better, faster, stronger. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 7263–7271).

Wang, C. Y., Liao, H. Y. M., Wu, Y. H., Chen, P. Y., Hsieh, J. W., & Yeh, I. H. (2020). CSPNet: A new backbone that can enhance learning capability of CNN. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops* (pp. 390–391).

Woebbecke, D. M., Meyer, G. E., Von Bargen, K., & Mortensen, D. A. (1995). Shape features for identifying young weeds using image analysis. *Transactions of the ASAE*, *38*(1), 271–281.

Yang, L., & Zhang, D. (2020). Real-time visual localization of the picking points for a ridge-planting strawberry harvesting robot. *IEEE Access*, *8*, 116556–116568.

Yu, Y., Zhang, K., Liu, H., Yang, L., & Zhang, D. (2020). Real-time visual localization of the picking points for a ridge-

planting strawberry harvesting robot. *IEEE Access*, *8*, 116556–116568.

# Instructions (delete all instructions, blue text, on all pages )

This Microsoft Word document includes particular Word styles, provided by ASABE, that are required for internet search indexing of your paper. The same Word styles provide a professional appearance. If you type or paste your material into the correct places (in the index boxes and at "Click here...") the styles will be correct automatically.

You can see the styles in Word on the Home tab in the Styles section. To see styles next to your text, go to File / Options / Advanced. In the Display section, in the box beside "Style area pane width in Draft and Outline views," enter a value of "1" and select "OK." On the View tab, select "Draft." Styles will appear to the left.

***Do not change the ASABE styles for the indexing boxes, title, authors, address, conference information, abstract, keywords, and references.*** The default, required, style names for these elements are Authors, Affiliation, Meeting Info, Paper Number, Title, Author(s), Address, Conf Name, Conf Sponsor, Conf Location, Conf Date, Other Pres, Abstract, Keywords, Ref Title, and Ref Listing.

ASABE styles for the body of your paper are Normal, Heading 1, Heading 2, Heading 3, Equation, Figure, Figure Caption, List Bullet, List Number, Table Caption, Table Contents, and Footnote. You may use other styles in the body of your paper.

Click these links for additional information about Word Styles and for other information for authors.

# Main Body (main headings use this "Heading 1" style)

Start a new paragraph with a single click of the Enter key, without a tab. Paragraphs will be indented.

Use styles for normal text, headings, figures, tables, captions, lists, etc. You may also use *italics*, **bold**, underlines, superscripts and subscripts. Generally use the Times New Roman font. For Greek letters and special symbols, use the Symbol font where possible. Avoid unusual symbols.

Type your text, highlight it, and select the appropriate style from the Styles tab. The text will change to the proper format when you apply the style.

Citations in the body of the text use the name, date system. For example, "Brown (2016) stated that... while others (Smith, 2011; Smith and Jones, 2013; Jones et al., 2014) found that....".

## Secondary Headings (this text is in the "Heading 2" style)

### Safety Emphasis (this text is in the "Heading 3" style)

You are urged to discuss the effects of your research, concept, design, technique, material, etc., on personal safety, if applicable. In what ways did you consider safety in your project? How will your work improve safety? What precautions do you plan or recommend for eliminating the adverse effects?

## Equations

Use the Word **"Equation"** style. Plain text may be used for a simple equation. MathType is preferred for equations, but you may use the built-in Word equation editor. Put the equation reference number outside the equation editor box. Tabs are set up to center the equation and to place the equation number at the right margin.

This is a plain text equation using the **"Equation"** style, with tabs before and after:

$$E = mc^2 \tag{1}$$

where
$E$ = kinetic energy
$m$ = mass
$c$ = the speed of light

## Figures (graphics, photos, charts, etc.)

Figures generally follow the paragraph where they are first mentioned. Use the Word **"Figure"** style for the image. Have a caption under each figure using the Word **"Figure Caption"** style. See example below.

The Word drawing canvas is best avoided. Only use it when absolutely necessary, namely, to constrain floating pieces,

such as arrows, within a figure.

For digital camera images (JPEG), use the medium or large file setting, not the small file (low quality) setting. For scans and other images, use 600 dpi for black and white line art or 300 dpi for color or grayscale. Higher resolution will not increase the quality of the published image.

Color figures will display in color in the web version but are printed in grayscale. Please print your color figures as grayscale and verify that they can be interpreted correctly, since the loss of color may make lines and gradients indistinguishable.

The font used inside a figure is different than the font used for the text of this paper or the figure caption. *Use a sans serif font, such as Arial, for all lettering inside figures to provide better clarity.* After making the image the size you want, the font within the figure should be 6 to 8 points. The caption uses Times New Roman font.
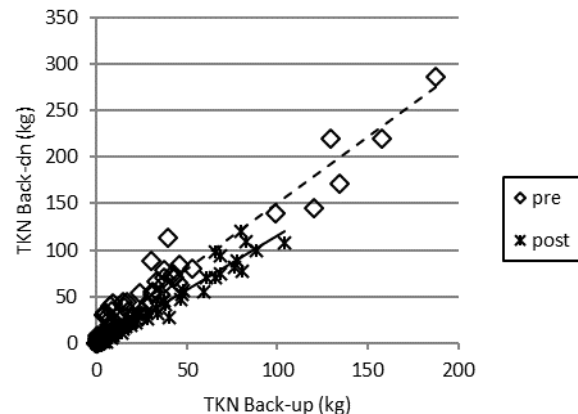


**Figure 1. Use the "Figure Caption" style for the caption below each figure. The figure caption should be separate from the graphics image. You may put your figures in tables to aid layout. Use Times New Roman font for this figure caption.**

## Tables

Tables generally follow the paragraph where they are first mentioned. Tables use the **"Table Contents"** style. The caption at the top of each table uses the **"Table Caption"** style.

**Table 1. Use the "Table Caption" style for the text. Material in the table uses the "Table Contents" style. Use the Word table tools or copy tables from Excel.**

| Material in the table uses the "Table Contents" style. | Internal line weight 0.5 point. The lines at the top and bottom of the table are 1 point[a] |
|---|---|
| | There is no line below the footnotes. |

[a] Footnotes use superscripted letters in brackets. Order them left to right, then the next row left to right, etc.

## Lists

You may use the **"List Bullet"** or **"List Number"** styles for your lists. Type the list, pressing Enter between items. Select all the listed items and apply the style. If Word forces text into the list against your wishes, press Backspace or select the text and make it **"Normal"** style.

- This uses the **"List Bullet"** style. Use bullets for lists unless numbering is necessary.
1. This uses the **"List Number"** style. Use a numbered list only when the list represents a sequence, such as the steps in a procedure.

# Conclusion (uses "Heading 1" style)

The Conclusion or Summary section restates the major findings and suggests further research. It is the last main heading before the references.

## Acknowledgements (uses "Heading 2" style)

Put any acknowledgements, such as thanks to contributing individuals or organizations, here.

# References (uses "Heading 1" style)

**[Click here to enter references]**

This section (for all references) uses the *APA 6th* style. Reference citations in the text use the name, date system. See those examples in the "Main Body" section above.

**NEW—In this section, we encourage the use of** <u>EndNote</u> **(preferred) or the** <u>Microsoft Word References tool in Word 2007 and later</u>**. See this** <u>video</u> **or** <u>these instructions</u>**. Choose** *APA 6th* **style and the** <u>ISO abbreviations for journal names</u> **(LTWA).** <u>JabRef</u> **(free download) can** <u>automate</u> **journal abbreviations for the ISO standard. JabRef can** <u>export references</u> **to "MS Office 2007 (*.xml)" that can easily be imported into Word References.**

**If you use the Word References tool or EndNote and select** *APA 6th***, all of the details of format described below will be done for you automatically when you create a bibliography.** Make it the **"Ref Listing"** Word style to create the indents.

**For those of you NOT using EndNote or Word References tool,** please compose your reference entries following the examples below, which follow the *APA 6th* style, and put the references in alphabetical order. (Word will do this for you under the Paragraph function).

Make the reference list the **"Ref Listing"** style.

List author names with last name first, then initials. List up to seven authors per author group; for sources with more than seven authors list the first six names, then an ellipsis (periods and spaces, . . . ), then the name of the last author, as in the first example.

For all titles, capitalize the only the first word, the first word after a colon or dash, and proper nouns. Titles of books and journals are italicized; other titles (book chapters, journal articles, papers from a meeting, reports, standards) are not italicized.

Do not use a period after a URL.

**Examples:**

**Journal Article**

Firstauthor, A. B., Second-Author, E., Thirdauthor, F. G., Fourth, H. I., Fifthauthor, J., Sixthauthor, K.L., . . . Lastauthor, Z. (2017). Title of journal article: Capitalize after colon. *Appl. Eng. Agric., 578*(12), 5-10. http://doi.org/10.1111/11111

**Book**

Lastname, A. B., & Jones Jr., C. D. (2014). *Book title.* Location of publisher: Publisher.

Surname, X. B., Author, C. D., & Jones Jr., E. (2017). *Book title* (2nd ed., Vol. 3). St. Joseph, MI: ASABE.

**Section of a Book**

Author, A., & Secondauthor, B. C. (2017). Section or chapter title. In *Book title* (pp. 17-34). Washington, DC: USEPA. Retrieved from www.epa.abcdefg.gov

Chapterauthor, A. (1987). Section or chapter title. In B. Bookeditor & C. Bookbinder (Eds.), *Book title* (2nd ed., Vol. 3, pp. 17-34). Rome, Italy: FAO.

**Paper from a Meeting, Conference Proceedings** Include the name and location of the publisher, but not the location or dates where the meeting was held. Abbreviate Conf., Int., Proc., Symp.

Author, A. B. (2018). Title of paper. In B. Editor (Ed.), *Proc. 10th Int. Conf. Agricultural Engineering. 2*, pp. 655-766. Washington, DC: USDA-ARS. http://doi.org/10.12/4x.57

Author, A. B., & Name, C. D. (2014). Title of paper [in Chinese]. *Proc. 10th Symp. Agricultural Engineering.* Publisher location: Publisher.

Author, A. B., & Name, C. D. (2016). Title of paper. ASABE Paper No. 1601234. St. Joseph, MI: ASABE.

**Standard**

ASABE Standards. (2018). S358.2: Moisture measurement—Forages. St. Joseph, MI: ASABE.

ASTM. (2014) 12343: Standard name. West Conshohocken, PA: ASTM Int.

**Dissertation or Thesis**

Author, A. (2014). Title of dissertation. PhD diss.[or MS thesis.] City, state or nation: University Name, Department of Engineering.

**Miscellaneous** If no author listed, use the name or abbreviation of the organization.

ABCD. (2014).Title. Association of BioCropsDiversity. Retrieved from http://bcd.org/report.pdf

Author, A. B. (2018). Patent title. U.S. Patent No. 123,456.

SAS. (1990). SAS User's Guide: Statistics. Ver. 6a. Cary, NC: SAS Institute.

USDA-NASS. (2004). Report title. Bulletin 1234. Washington, DC: USDA-NASS. Retrieved from www.usda.nass.gov/x1234.pdf

**Unpublished Material** Do not list such material in the References section because it is not available to the reader. Put useful information in the text of your manuscript, e.g.,

". . . this was rare (Charles Brown, USDA-ARS, personal communication, 23 May 2018).

# Appendix or Nomenclature (optional)