

# Manuscript Template

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

## FRONT MATTER

### Title

Full titles : Enhanced field-based detection of potato blight in complex backgrounds using deep learning

Short titles : Mask-RCNN based potato blight detection

### Authors

Joe Johnson<sup>1</sup>, Geetanjali Sharma<sup>1</sup>, Srikant Srinivasan<sup>1\*</sup>, Shyam Kumar Masakapalli<sup>2</sup>, Sanjeev Sharma<sup>3</sup>, Jagdev Sharma<sup>3</sup>, Vijay Kumar Dua<sup>3</sup>

### Affiliations

<sup>1</sup>School of Computing & Electrical Engineering, Indian Institute of Technology Mandi, Kamand, H.P., India

<sup>2</sup>BioX Center, School of Basic Sciences, Indian Institute of Technology Mandi, Kamand, H.P., India

<sup>3</sup>ICAR-Central Potato Research Institute, Shimla, H.P., India

\*Corresponding author. Email: srikant\_srinivasan@iitmandi.ac.in

### Abstract

Rapid and automated identification of blight disease in potato will help farmers to apply timely remedies to protect their produce. Manual detection of blight disease can be cumbersome and may require trained experts. To overcome these issues, we present an automated system using the Mask Region-based convolutional neural network (Mask R-CNN) architecture, with Residual Network as the backbone network for detecting blight disease patches on potato leaves in field conditions. The approach uses transfer learning, which can generate good results even with small datasets. The model was trained on a dataset of 1423 images of potato leaves obtained from fields in different geographical locations and at different times of the day. The images were manually annotated to create over 6200 labelled patches covering diseased and healthy portions of the leaf. The Mask R-CNN model was able to correctly differentiate between the diseased patch on the potato leaf and the similar looking background soil patches, which can confound the outcome of binary classification. To improve the detection performance, the original RGB dataset was then converted to HSL, HSV, LAB, XYZ, and YCrCb color spaces. A separate model was created for each color space and tested on 417 field-based test images. This yielded 81.4% mean average precision on the LAB model and 56.9% mean average recall on the HSL model, slightly outperforming the original RGB color space model. Manual analysis of the detection performance indicates an overall precision of 98% on leaf images in a field environment containing complex backgrounds.

43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
  
85  
86  
87  
88  
89

## MAIN TEXT

### 1. Introduction

Early and late blight diseases are a common occurrence across regions where potato (*Solanum tuberosum* L.) is cultivated. Blight is a common foliage disease of potato that starts as uneven light green lesions near the tip and the margins of the leaf and then spreads into large brown to purplish-black necrotic patches as reported by Arora et al. [1]. Blight causes premature defoliation and eventually incites tuber rot of potato. As noted by Haverkort et al. [2] unchecked blight could destroy the entire crop within a week under conducive conditions. Thus, blight in potato could bring disastrous consequences, particularly to farmers with marginal landholding who grow potato as cash crops [3].

In most developing countries, detection and identification of blight is performed manually by trained personnel scouting the field and inspecting potato foliage. This process is tedious, and in some cases impractical, due to the unavailability of a disease expert in remote regions [4]. On the other hand the recent advances in image processing for rapid and automated disease identification using images of plant leaves [5-7] can make the process far more efficient and timely. In the recent past a system has been proposed to identify the severity of potato late blight disease from field images using Fuzzy C-mean clustering [8] but with few images. Using 300 images as a training set, another work [3] has attempted potato disease detection using segmentation and multi-class support vector machine. These datasets do not incorporate time-varying illumination and are usually taken at a fixed time corresponding to the best illumination. Usually methods developed using small datasets do not perform well in field environments due to the large variations in illumination, focus, resolution, underlying feature size and presence of occluding objects in the images.

More recently, the task of classification and detection in images has been dominated by various flavors of neural networks (NNs), especially with the advent of deep NNs [9-14]. It has been well-accepted that deep learning models perform quite well in image classification and detection compared to traditional image-processing algorithms [15]. The process of trial and error for fine-tuning traditional image processing models to obtain the representational features of objects becomes rapidly complicated as the number of classes increases. On the other hand, a neural network learns complicated underlying patterns specific to a certain class of object without any manual intervention. A classification model using convolutional neural network (CNN) for distinguishing 58 classes of healthy and diseased plant dataset was developed by Ferentinos [16]. Arsenovic et al. [17] have improved the plant disease classification by increasing the training dataset, which has images of leaves in field conditions. Deep NNs have the potential to quickly detect an object from a complex image, which makes them suitable for smart-phone applications [7]. However, the training process in deep NNs is computationally expensive where the network parameters are iteratively fine-tuned to improve the mapping between a set of training input images and the desired outputs [9]. Therefore, such methods have become popular only with the concomitant advances in graphics processing hardware.

In the context of an image comprising a potato leaf amidst a complex background, the classification process has a binary outcome, i.e., it determines if the overall image reflects disease or not. Detection, on the other hand, goes one step further and demarcates the specific patch or patches on the leaf that contain the signature of blight. Region-based deep CNN (R-CNN) [18] is an object detection method that is trained to propose regions by

90 exhaustively searching the image after it has been transformed through several convolution  
91 layers. For the purpose of object detection, architectures like YOLO [19], SSD [20], Faster  
92 R-CNN [21] and Mask R-CNN [22] are recent methods, with Mask R-CNNs giving a better  
93 overall performance. For the R-CNN architectures, the Residual network with 50 layers  
94 (ResNet-50) is usually used as a backbone. Other applications of CNN in agriculture include  
95 Zhang et al. [23] who have used global pooling dilated CNN for better segmentation and  
96 classification of cucumber leaf disease, while CNN-based regression has been used to  
97 estimate soybean leaf defoliation with the aid of real and synthetic images [24].

98 The various transformations that an image undergoes as it traverses a deep CNN can  
99 sometimes be understood by visualizing the output of individual convolutional layers. The  
100 output is termed as the feature map or activation map and can be visually correlated to the  
101 input image. Each convolutional layer is a set of functions that applies some transformation  
102 to the image, behaving as a filter. The feature map aids in relating the learned filter with the  
103 performance of the model and using the learned filter to improve the performance as  
104 discussed in [25]. Lee et al. [26] have reported such studies in plants where the different  
105 orders of venation provide better representative features than the outline shape of a leaf  
106 when considering the hierarchical transformation of features from lower-level to higher-  
107 level abstraction for species classes.

108 In addition to the choice of appropriate NN architectures, preprocessing the image data can  
109 contribute towards obtaining better detection or classification. For example, it has been  
110 observed that a color spectrum provides better results than grayscale for object detection by  
111 deep learning models [7]. A color space or color model is mathematical transformation to  
112 project a set of primary colors to a different range of colors [27]. An investigation of the  
113 influence of different color spaces to improve the deep-learning model performance has  
114 been conducted for the traffic light detection system [28]. A comparative study for different  
115 color space using deep learning-based automatic segmentation system has been discussed  
116 in [29]. Disease region segmentation of paddy crop using Mask-RCNN on different color  
117 space images are analyzed in [30]. Robustness and accuracy of the segmentation of foliar  
118 disease spots images using region growth and comprehensive color features has been  
119 explored in [31].

120 The objective of this work is to develop a Mask R-CNN based model to detect the blight  
121 symptoms on an infected potato leaf, which can eventually be deployed on a cell phone.  
122 Mask R-CNN [22] is chosen because it utilizes a Feature pyramid network (FPN), allowing  
123 it to grasp semantically relevant features at different resolution scales. The region proposal  
124 network (RPN) scans the entire top-bottom pathway of the FPN for feature maps containing  
125 required objects and proposes regions of interest (ROI). This enables prediction of relevant  
126 classes, bounding boxes, and mask for the region or patch. These methods of Mask R-CNN  
127 force different layers in neural network to learn features across multiple scales, making it  
128 robust to several environmental variations in the image. The model learn features from  
129 visual characteristics such as the shape, color, texture, and venation of potato leaf and blight  
130 disease for different training data.

131 The emphasis on detection rather than classification is because simple classification into  
132 healthy or unhealthy categories can be misleading due to misclassification of soil patches in  
133 the background as disease. To improve blight detection, we also investigate preprocessing  
134 the data to include different color space images. Fig. 1 conveys the overview of the method  
135 proposed in this work. We have converted the RGB color space dataset to five other color

spaces, namely, HSL, HSV, LAB, XYZ, and YCrCb and created a separate Mask-RCNN model for each color space. The model uses transfer learning or stored knowledge of a pre-trained Mask R-CNN model on the Microsoft Common Objects in Context (MS COCO) dataset [32] as the initial condition for the training process. The performance of the networks across the different color spaces are compared in their ability to automatically detect infected potato leaves and disease patches in complex field images.

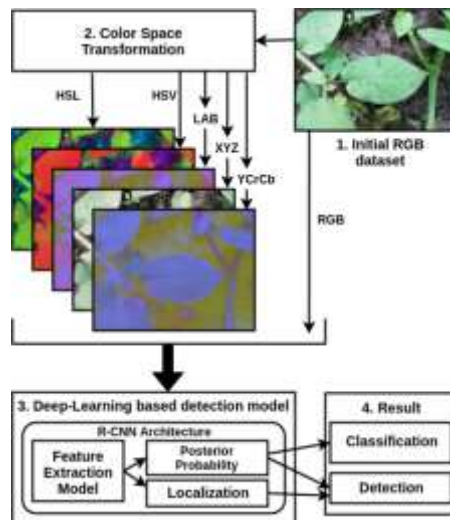


Figure 1: Overview of the deep-learning based potato disease classification and detection method in this work.

## 2. Materials and Methods

### 2.1. Data acquisition

The choice of data used for training a CNN has a very strong impact of the effectiveness of the model in different situations. Factors such as the characteristics of the imaging sensor, the imaging protocol followed, illumination variation due to time of the day, shadows due to nearby objects, occlusion and complex background information, all need to be carefully considered to create a model that can be successfully applied to field-based imaging. In order to maximize the diversity of training data a set of 1840 field-based images of potato leaves was acquired for this work across different states in India by field personnel deputed under the FarmerZone project [33].

The dataset comprises images of healthy potato leaves as well as leaves affected by both early and late blights. As one of the objectives was to develop a model that would be accessible to a larger group of small-scale farmers, it was determined that the choice of imaging sensors should include low-end cellular phones. Therefore, the potato leaf dataset contains images of resolutions of 3072 x 4096 pixels (552 images), 3120 x 4160 pixels (922 images), and 2448 x 3264 pixels (366 images) due to inherent differences in the sensors of the different smart-phones used for data collection.

The images are heterogeneous, having been taken from different locations within the field at different times of the day, typically between 11 am and 2 pm. Each image can contain several leaves, soil, and weeds in the background apart from the primary infected/healthy potato leaf. This variation aids the generalization of the deep learning model. All images were captured in natural light with the camera flash always turned

off, and without any additional optical or digital zoom. Sample images of healthy leaves and leaves affected with blight are shown in Fig. 2.

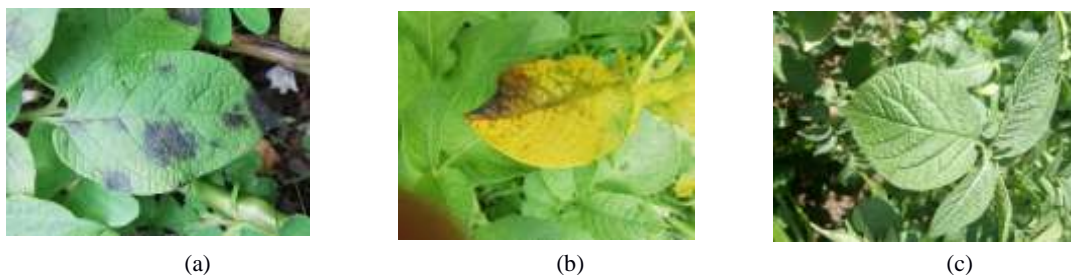


Figure 2. The complex background dataset used in this study. (a) Multiple disease patches on a single leaf. (b) Single disease patch on a discolored leaf. (c) Healthy leaves

## 2.2. Data curation

The potato leaf images obtained using smart-phones are in the RGB format, which is similar to the human perception of the light spectrum as a combination of the primary colors – red, green, and blue [34]. While there is potential for improved image segmentation using other color spaces, there is no general opinion on the best choice of color space for image segmentation. Therefore, all the RGB images were converted to five color spaces (HSV, HSL, XYZ, LAB, and YCrCb) using Open Computer Vision Library [35], creating an additional 5 datasets. In the RGB dataset, one or more blight spots on each potato leaf in the foreground are manually demarcated into patches for creating the ground-truth dataset. The process of demarcating or segmenting the images was carried out by three personnel, two non-experts under the guidance of an agricultural expert. The ground-truth values and labels are kept same for all the images in different color space datasets. To reduce the annotator’s bias and variance [36] during ground-truth annotation, the following steps were taken:

- i) The expert first demonstrated the protocol for segmentation of patches on foreground, the edges to be considered and how tightly the polygon should be drawn.
- ii) For 50 randomly selected images, the expert and the non-experts all annotated according to the prescribed procedure. The value of Cohen’s Kappa [37] found across the three annotators was 0.92 and level of agreement was found to be very good. Thereafter, 5825 blight patches, 356 infected leaf patches, and 89 healthy leaf patches were created from the 1840 input images. Table. 1 provides the details of the total dataset and its annotation count. To create and validate the disease detection model, the dataset of each colorspace was further split into 2 sets containing approximately 80% and 20% data, respectively.

Table 1: Description of manually annotated patches for different features in the dataset

Data	Train	Test	Total
<b>No. of Images</b>	<b>1423</b>	<b>417</b>	<b>1840</b>
No. of blight patches	4673	1152	5825
No. of Infected leaf	1423	356	1779
No. of Healthy leaf	122	89	211

## 2.3. Mask R-CNN based detection model

The detailed block diagram of Mask R-CNN used in this work is shown in Fig. 3. Mask R-CNN is an extension of Faster R-CNN [38], with an additional forking to a prediction segmentation mask on each RoI, in parallel with the already available branch for

classification and bounding box regression. In this work, further tuning of the original Mask RCNN include the use of ResNet-50 as backbone architecture with RPN anchor scales set to 32, 64, 128, 256, and 512 and the anchor aspect ratios set to 1:2, 1:1, and 2:1. This follows from manual observation of the training dataset, which shows that the various demarcated patches vary in this selected range of pixel values and aspect ratios.

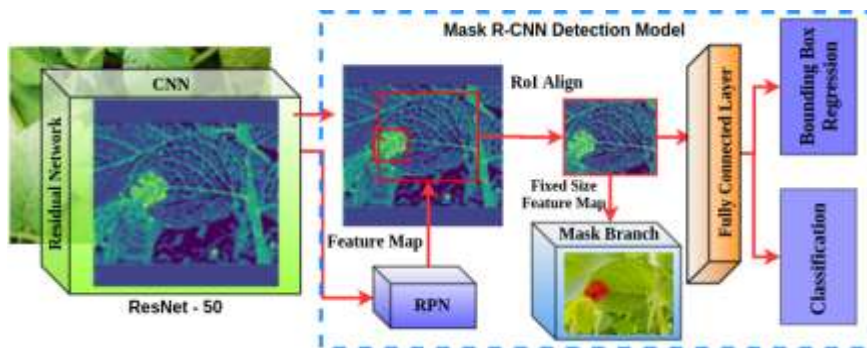


Figure 3: Block diagram of a model architecture for the implemented binary classification and Mask R-CNN models.

Regarding the choice of ResNet-50 as the backbone, it may be noted that deep CNN is prone to problems like vanishing gradients and the curse of dimensionality [14], with an increase in the number of layers. To avoid this degradation problem for a deeper network, skip connections (identity connections), or residual connections are used. The residual connection is a ‘shortcut’ module whereby, the weight/convolutional layers are skipped and the input is added through an identity function before the final ReLU activation function. It is observed that during backpropagation, larger gradients are available for initial layers leading to faster learning because of skip or residual connection. ResNet-50 has 50 layers arranged in five stages with a total of sixteen residual blocks. In each residual block, the convolutional layer is followed by a batch normalization layer and a ReLU activation function. The ResNet-50 model generates 256, 512, 1024 and 2048 feature maps from the second, third, fourth and fifth stages respectively.

Each color space dataset is used for training separate Mask R-CNN detection model. In a preprocessing step, the input images are down-sampled to 1024 x 1024 pixels. For each color space model the mean value of each channel of the respective color space, calculated separately from the training dataset, is set in the configuration file of the program [38]. Pre-trained weights of the MS COCO dataset have been used for the initial training of the model as attempts to train from scratch did not yield significant detection even after 70<sup>th</sup> epoch for all of the color space datasets, probably due to the small dataset. On the other hand, the application of transfer learning towards classification of potato leaf disease was shown in [39,40]. To optimize the network weights, the stochastic gradient descent optimizer with momentum fixed at 0.9 was used. A fixed learning rate of 1e-4 was set for optimum learning. The maximum number of epochs was set to 100 and iterations per epoch were set to 712 corresponding to a batch size of two images per iteration.

#### 2.4. Computing resources utilized

The training and testing of the model was performed on a CentOS 7 Linux workstation equipped with one Intel Xeon Processor CPU (96 GB RAM), accelerated by one Nvidia GeForce GTX 1080 Ti GPU (11GB Memory). The model is implemented in the Keras 2.2.4 deep learning open-source framework with the Tensorflow-GPU 1.8.0 backend using

242 Python 3.6. The detection model on each color space took an average of 25 hours for  
243 training.

244 For creating the ground truth dataset VGG Image Annotator (VIA) [41], a standalone  
245 software, was used for the manual annotation of the blight and leaf patches in the image. It  
246 allows rectangular and polygonal shaped area to be annotated, which is useful for training  
247 Mask R-CNN.

## 248 2.5. Model evaluation metrics

249 In computer vision, standard metrics like precision and recall are used for performance  
250 evaluation of binary classification [42]. This is obtained from a confusion matrix that  
251 summarizes the performance of a classifier for a given test dataset. The four components of  
252 the 2x2 confusion matrix for any binary classifier are True positive (TP), True negative  
253 (TN), False Positive (FP) and False negative (FN). The correct classification of an image  
254 containing disease would count as a TP, while an incorrect classification as a healthy image  
255 would count as a FP. The performance of the classifier is then obtained by:

$$256 \quad \textit{Precision} = \frac{TP}{TP+FP} \quad \textit{Recall} = \frac{TP}{TP+FN}$$

257 For the performance assessment of the object detection model, both the correct  
258 classification and the precise location of the disease patch in the image should be taken into  
259 account. To do so, concepts such as intersection-over-union (IoU) and the average precision  
260 (AP) was introduced in the Pascal VOC challenge [43]. The IoU metric determines the  
261 correctness of the patch detection by taking into account how closely the predicted instance  
262 (PI) fits the ground truth instance (G). IoU is the measure of overlap between G and PI  
263 boundaries given by

$$264 \quad \textit{IoU}(G, PI) = \frac{G \cap PI}{G \cup PI}$$

265 The IoU threshold is taken to be 0.5 as a common practice, whereby, if the IoU value of  
266 detection is greater than 0.5 then the PI is considered as a TP, else it is taken as a FP. This  
267 is illustrated using a sample test image shown in Fig. 4, where green color masks and  
268 bounding boxes represent the human-annotated ground truth while red color masks and  
269 bounding boxes represent the predictions by the detection model. For the sample image  
270 shown in Fig. 4, the confidence score and IoU for the infected leaf are 100% and 93 %  
271 respectively.

272 In addition to the boundaries of the PI, the algorithm also provides a confidence level for  
273 the PI. The AP is a metric that incorporates the confidence level of prediction and IoU into  
274 the calculation of precision using area under precision-recall curve. Mean average precision  
275 (mAP) is mean of AP across the different categories or classes, which are detected, and  
276 summarizes the performance of a detection model.



Figure 4: An example to describe IoU on the object with and without annotation.

### 3. Results and Analysis

#### 3.1. Disease detection

The performance of the disease detection model, when tested on the ground truth potato leaf dataset, is calculated according to the metrics defined in section 2.5. A separate model is created for each color space. Even within each color space, there are two types of Mask R-CNN models:

(i) Two-class model: this involves the detection of only potato blight patches, while the rest of the image is considered as background. This kind of demarcation is a natural first step where it is expected to detect only blight disease patches from the input image. However, once the model was trained over the entire dataset, it was found that several blight patches were not detected and that a few soil patches were misclassified as blight. A sample test image from the RGB dataset [Fig. 5(a)] contains nine disease patches spread across three different leaves. Fig. 5(b) shows that the two-class model has detected only two disease patches out of nine clearly distinguishable disease patches.



Figure 5: (a) Sample test RGB image and (b) output image for trained two-class Mask R-CNN model.

(ii) Four-class model: As a means to improve the performance of the detection model, a second experiment was performed in which the Mask R-CNN model was trained to detect 4 classes: blight disease patches, infected leaves, and healthy leaves, in addition to the background [Fig. 6, 7].



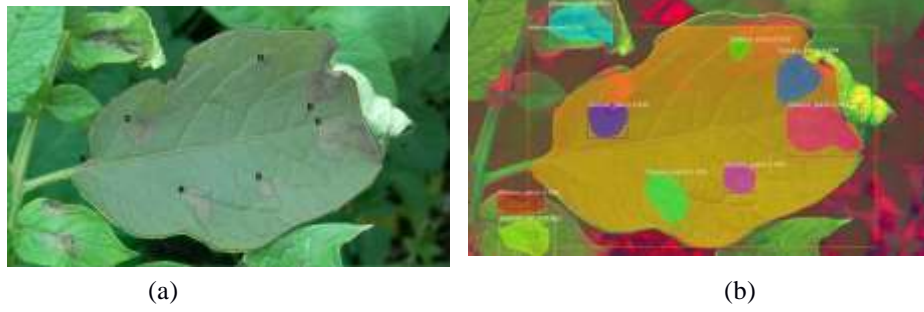


Figure 6: Sample RGB image (a) Human annotated foreground regions (b) Patches inferred by 4-class HSL model.

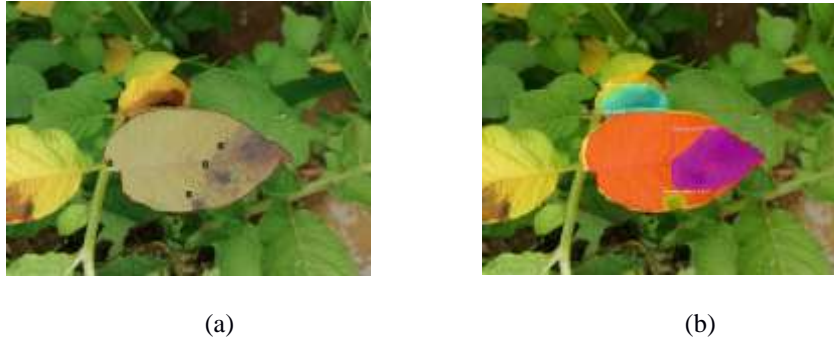


Figure 7: Sample RGB image (a) with annotated foreground regions (b) Patches inferred by 4-class RGB model.

For both models, the ground truth criteria were kept uniform for all the images. The aim of this second model was to increase blight disease patch detection and reduce the FP due to misclassification of soil as disease, by the inclusion of a post-processing step that checks for the intersection of the disease patch with the leaf patch. Nevertheless, it was seen that the performance of the four-class model was superior to that of the two-class model even without any additional postprocessing. The performance scores for the 2-class and the 4-class detection models are compared in Table. 3 with respect to different color spaces.

Table 3: Performance of various color spaces compared to RGB for potato disease detection

Color space	2-class detection model		4-class detection model		Avg. Inference time/image (sec)
	mAP (%)	mAR (%)	mAP (%)	mAR (%)	
RGB	77.4	53.9	80.9	55.5	1.77
XYZ	77.4	54.2	70.8	52.9	1.68
HSL	78.2	55.2	81.3	<b>56.9</b>	1.67
HSV	75.5	53.9	75.5	55.4	1.67
LAB	<b>80.1</b>	<b>55.6</b>	<b>81.4</b>	56.4	1.68
YCrCb	76.1	54.3	79.1	56.4	1.70

Performance score calculated for IoU = 0.5.

Among the two-class Mask R-CNN models for different color spaces, LAB color space has the best mAP (80.1%) and mAR (55.6%) values. The 4-class detection model shows a slightly improved mAP (LAB) and mAR (HSL) performance metrics of 81.4% and 56.9% respectively. It was observed that HSL, LAB and YCrCb color space models could perform

325 better than RGB color space model overall and specifically for disease patch detection.  
 326 Inference times are the least for the detection model trained on HSL and HSV color spaces.  
 327 These results are in line with the latest leaderboard on the COCO website [44], which  
 328 publishes the performance results for different models on the COCO dataset with 91  
 329 categories. The highest mAP (IoU=0.5) is 60.6 % for a model trained on the dataset of  
 330 broccoli category (closest to potato leaves).  
 331

332 Further investigation into the performance of different four-class detection models by  
 333 manually comparing the test image data to the model outputs shows that the performance  
 334 of the disease detection model appears far better than the mAP and mAR values reported in  
 335 Table 3. This surprising outcome can be understood if we delve into the ground truth  
 336 labelling procedures. The images taken from the field have many complex regions due to  
 337 fuzziness of image, partially occluded disease patches and disease patches on the stem.  
 338 Many disease patches that fall into these categories were not annotated while creating the  
 339 ground truth dataset. Also for the human annotator, there is often no clear distinction  
 340 between the foreground and background features, whether for disease patches or the leaves.  
 341 The human annotator, for example, has annotated (shaded region shown in Fig 6.a and 7.a)  
 342 only clearly distinct features of disease or leaf patches. However, our trained models have  
 343 correctly predicted several unlabelled disease patches in the background as disease [Figs 6b  
 344 and 7b]. Since these ‘vague’ disease patches have not been labeled in the ground truth  
 345 dataset, they end up lowering the mAP and mAR parameters, despite the correct  
 346 classification by the model. Hence, the ground truth annotation might need to be more  
 347 inclusive of disease patches to better represent the performance score.  
 348

### 349 3.2. Manual analysis of the detection model

350 Considering the challenges of ground truth labelling, we attempted a more realistic  
 351 quantification of disease patch detections by manually verifying the outputs of the 4-class  
 352 model (Table. 4). The correct disease patch predictions were categorized into two true  
 353 positive categories: TP1 reflects the correct detections that match the ground truth while  
 354 TP2 reflects correct disease detections that have not been annotated in the ground truth. The  
 355 same exercise is carried out for the infected leaf and healthy leaf classes also. Table. 4  
 356 summarizes the results of manually determining the detection performance on the test  
 357 dataset, for all color spaces. It can be inferred that among the six color space models, the  
 358 model trained on HSL color space has best disease detection with 464 (TP1) patches  
 359 detected. The LAB and YCrCb color space trained models have the best combined four-  
 360 class performance metrics of 98.6% (Combined precision) and 85.8% (Combined recall)  
 361 respectively. The YCrCb model shows maximum true disease detection (TP1+TP2) of 647  
 362 disease patches. The infected leaf patches were detected better by HSV color space model  
 363 with 341 true detections. It was observed that HSL, HSV, LAB, and YCrCb models  
 364 performed better than the RGB color space model for the detection of disease patch and  
 365 infected leaf. In all color space instances of the 4-class model, very few FP are observed  
 366 for disease patch class while most of the FP in the infected leaf class are misclassifications  
 367 of a healthy leaf.  
 368

Table 4: Manually obtained performance metrics of the four-class Mask R-CNN model

Color space	Disease patch				Infected leaf			Healthy leaf			Combined	
	TP1	TP2	FN	FP	TP	FN	FP	TP	FN	FP	P(%)	R(%)
RGB	375	166	176	7	329	22	10	93	15	2	98.1	81.9
XYZ	343	191	173	1	340	11	34	52	55	1	96.3	79.5
HSL	464	159	147	5	338	13	21	83	25	2	97.5	85.0
HSV	428	149	140	7	341	10	23	68	40	1	96.9	83.8


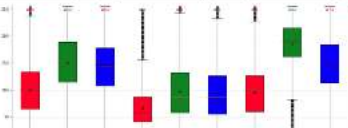
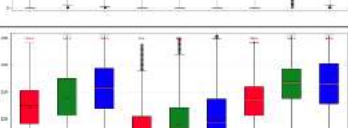
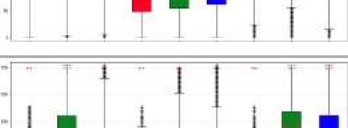
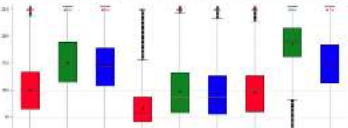
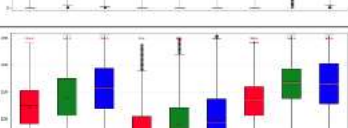
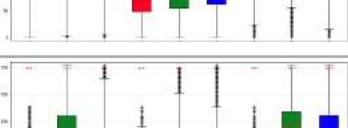
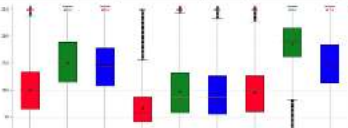
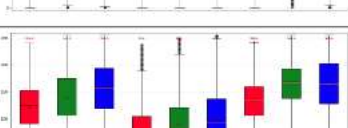
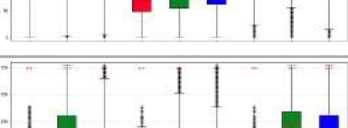

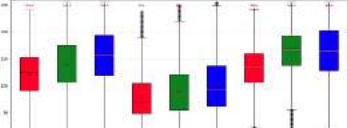
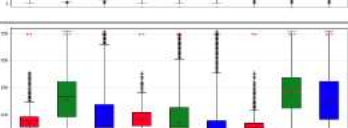
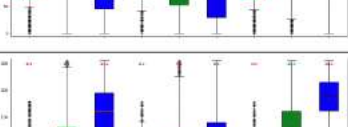
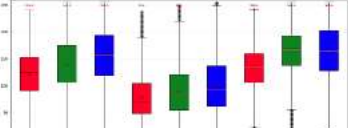
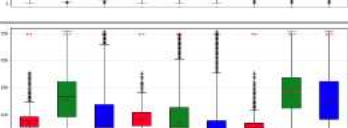
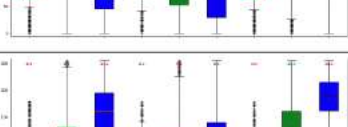
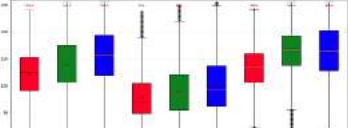
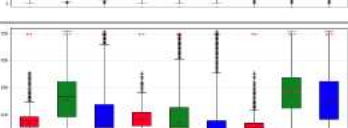
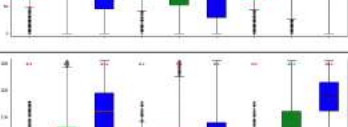
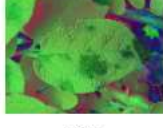
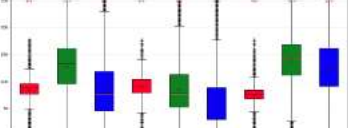
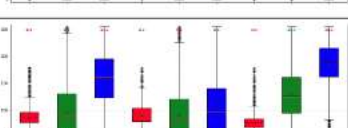
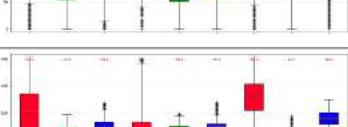
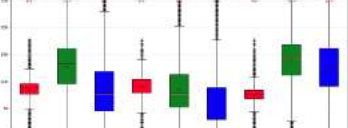
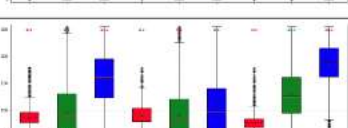
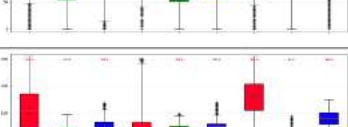
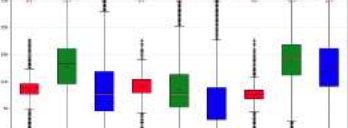
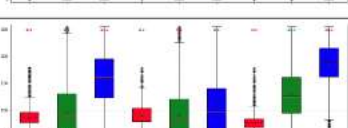
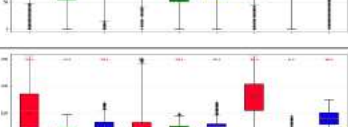
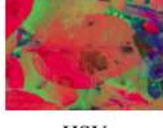
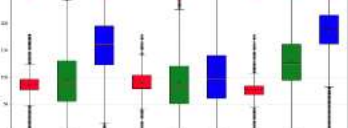
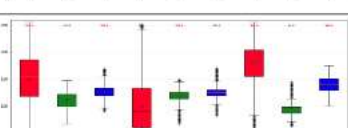
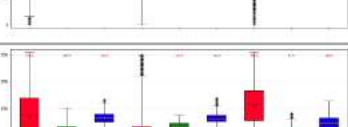
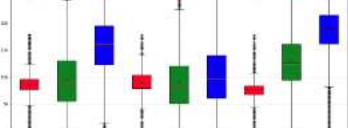
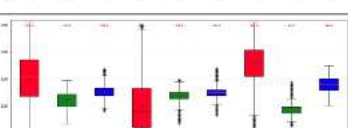
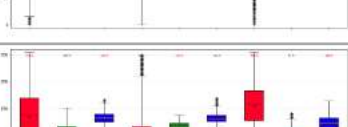
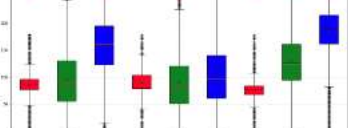
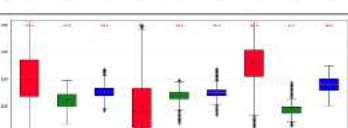
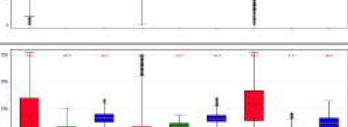
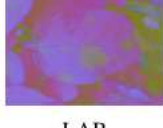
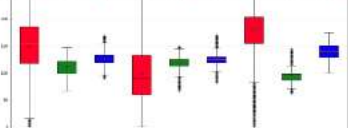
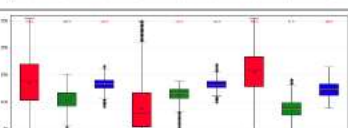
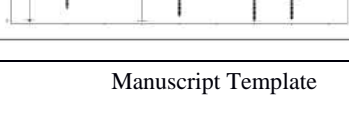
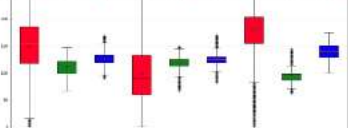
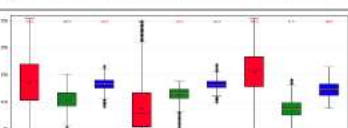
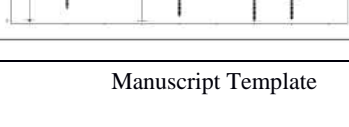
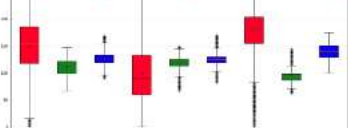
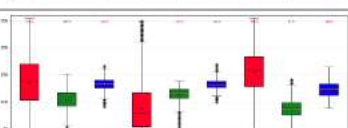
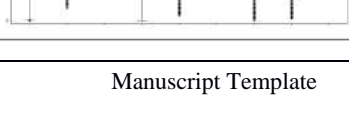
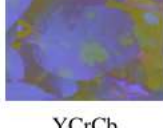
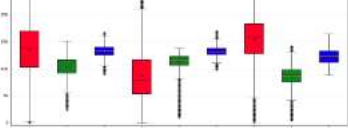
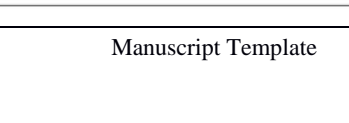
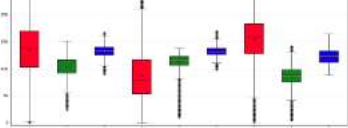
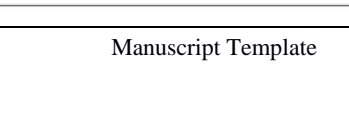
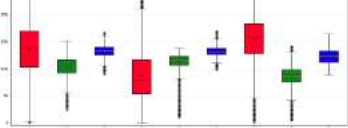
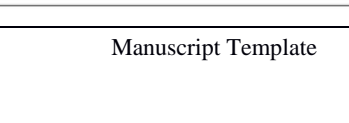
LAB	395	175	180	2	333	18	11	91	17	1	<b>98.6</b>	82.2
YCrCb	394	253	101	9	336	15	30	52	56	0	96.4	<b>85.8</b>

369

### 370 3.3. Analysis of the role of color spaces

371 Hadji et al. [45] have previously shown that the histogram of image intensities is used  
 372 broadly for recognition and retrieval in an image database. For a better understanding of the  
 373 effect of each color space on the potato leaf dataset, the histogram trends of various color  
 374 components in the image can be observed [46]. Table 5 shows a histogram analysis on thirty  
 375 randomly selected potato leaf images with blight symptoms. Each image was of 2448 x  
 376 3264 pixels and further divided into image patches of 200x200 pixels. All patches were  
 377 manually labeled into classes of blight disease, healthy leaf, soil, and background. The count  
 378 of patches for each class were as follows: the blight disease class contained 844 patches,  
 379 the healthy leaf class contained 2216 patches, the soil class contained 102 patches and  
 380 patches that did not meet the predetermined patch size and features were discarded. The  
 381 pixel intensity distribution for patches of disease (D), soil (S) and leaf (L) regions in the  
 382 form of a box plot with the mean ( $\mu$ ) and the standard deviation ( $\sigma$ ) for each channel of the  
 383 color space is shown.

384 Table 5: Sample image, box plot for pixel intensity distribution for three-channel among different color spaces with a  
 385 mean ( $\mu$ ), and standard deviation ( $\sigma$ ) for three features for the different color spaces. (Boxplot representation for channel  
 386 1, channel 2 and channel 3 are represented by red, green and blue color respectively)

Color space with Sample Image	Boxplot for pixel intensity distribution			Region	Channel			
	Disease (D)	Soil (S)	leaf (L)		1	2	3	
 RGB				L	$\mu$	95	188	146
				S	$\sigma$	45.6	36.9	46.8
				S	$\mu$	68	99	98
				D	$\sigma$	37.5	51.4	48.2
				D	$\mu$	101	151	144
					$\sigma$	46.9	49.4	48.8
 XYZ				L	$\mu$	130	165	165
				S	$\sigma$	35.6	36.2	47.7
				S	$\mu$	77	92	102
				D	$\sigma$	40.1	46.3	51.7
				D	$\mu$	120	140	155
					$\sigma$	41.9	45.9	52.0
 HSL				L	$\mu$	75	141	125
				S	$\sigma$	11.3	37.3	48.2
				S	$\mu$	91	82	65
				D	$\sigma$	19.4	42.8	49.0
				D	$\mu$	85	130	82
					$\sigma$	20.4	44.3	52.3
 HSV				L	$\mu$	75	127	188
				S	$\sigma$	11.3	37.3	48.2
				S	$\mu$	91	82	101
				D	$\sigma$	19.4	42.8	49.0
				D	$\mu$	85	99	159
					$\sigma$	20.4	44.3	52.3
 LAB				L	$\mu$	125	94	140
				S	$\sigma$	34.2	9.8	12.7
				S	$\mu$	100	123	126
				D	$\sigma$	49.7	13.5	11.4
				D	$\mu$	150	112	127
					$\sigma$	46.4	13.9	12.2
 YCrCb				L	$\mu$	156	86	124
				S	$\sigma$	36.5	18.4	14.0
				S	$\mu$	76	115	131
				D	$\sigma$	44.6	14.9	9.1
				D	$\mu$	140	103	133
					$\sigma$	45.0	18.1	11.4

387

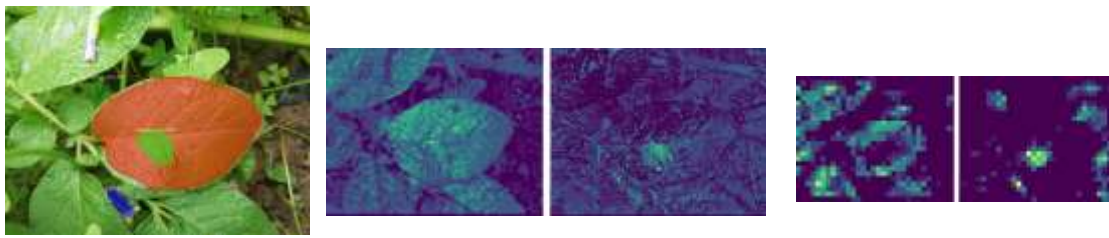
389 It is observed from Table. 5 that each of the channels of the RGB color space shows pixel  
390 distribution with a large spread that overlaps with the adjacent regions. This is due to  
391 varying illumination conditions across images, which equally affects the R, G, and B  
392 channels. Only for channel 2 is there some separation between the distribution of disease  
393 and leaf regions. This color information might be used by a deep learning model for  
394 classification. Similar to the RGB color space, the XYZ color space has a wide distribution  
395 of intensity values for all the components. Here, the channel 2 has better separation between  
396 soil and leaf regions.

397 Conversion to HSL from RGB color space restricts the range of certain components such  
398 as hue, which are illumination independent. Therefore, the hue component is expected to  
399 have a narrow range for all regions. Thus, it can be observed that although each of the HSL  
400 channels' distribution overlaps across all regions, the overlap is minimum for the hue  
401 channel. The many outliers might still make classification difficult using only hue  
402 information. Similar to HSL, the HSV color space has the same hue information. Compared  
403 to HSL, the HSV has more spread in pixel intensity distribution for channels 2 and 3. This  
404 might lead to reduced performance of the detection model. In the case of LAB, channel 1  
405 (lightness) varies according to lighting conditions. The component a\* and b\*, represented  
406 by green and blue boxes respectively, are the green-red component and a blue-yellow  
407 component of the image. From Table 5, it can be observed that a\* and b\* components have  
408 a narrow range of pixel intensity values. Here, the a\* component shows separation in values  
409 for leaf, soil and disease regions. Similarly, the b\* component has a very narrow  
410 overlapping area. This clear segregation in the a\* and b\* components could help in better  
411 classification and detection models. For the case of YCrCb, blue-difference and red-  
412 difference chroma have narrow spreads for different regions. The red box represents the  
413 luma component, green and blue boxes represent blue and red-difference chroma  
414 respectively. In channel 2, the soil and the leaf regions are distributed apart from each other  
415 while all other channels overlap in their distributions for the different regions. Thus, the  
416 histogram analysis helps to understand the color complexity of different patches/regions  
417 and it is observed that the color information will solely not lead to good segmentation  
418 between disease, soil and leaf regions. Higher level features like the texture of the disease  
419 region and leaf venation will need to be used by the deep learning model for segmentation,  
420 in addition to color.

421

### 422 3.4. Feature Map Observations

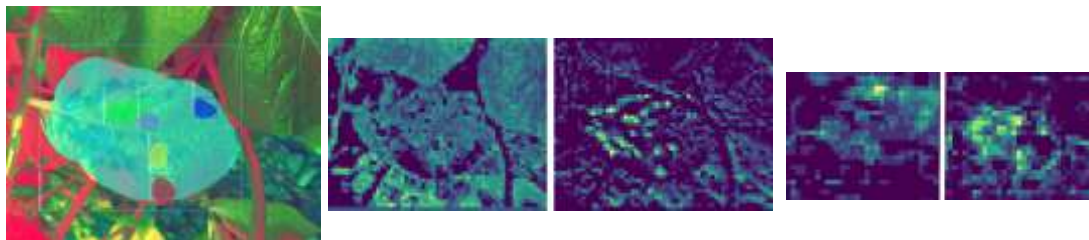
423 The characterization of blight disease, soil, and leaf regions by the CNN can be observed  
424 from the deconvolutional layers. The deconvolutional network maps the feature activity  
425 back to the input pixel space by using the same components of the convolutional layer  
426 (filtering, pooling) in the reverse order [25,26]. The feature maps of the different stages of  
427 ResNet, trained for four-class detection, include a number of relevant ones showing features  
428 of leaves and disease patches. Fig. 8.a shows the sample output image for the RGB color  
429 space detection model. Figs. 8.b and 8.c both show the visualization of the leaf feature and  
430 disease patch feature side-by-side for the 2c and the 5c layer respectively. From Fig. 8 it  
431 was inferred that the leaf features are well learned. The activation in the 5c layer shows that  
432 along with features of disease patch, features of soil patch have also been learned by the  
433 detection model. Overall, the 36<sup>th</sup>, 12<sup>th</sup>, 31<sup>th</sup> and 28<sup>th</sup> feature maps of the 2<sup>nd</sup> to the 5<sup>th</sup> layers  
434 were strongly related to disease, flower in the background, infected leaf, and healthy leaf  
435 respectively. The learning of leaf venation could be properly observed in the feature maps.



(a) RGB output image (b) 2c layer (c) 5c layer

Figure 8: (a) Sample output image, corresponding visualization of the leaf, and disease-relevant feature map for (b) 2c and (c) 5c layers of RGB color space detection model.

It is interesting to observe that for the model trained with HSL color space dataset, learning and extraction of the finer leaf and disease patch features were observed in the visualizations for a second to the fifth stage of the model. From the leaf feature maps, it could be inferred that a leaf's feature is better learned when it is in proper camera focus. The feature maps for the sample HSL output image is shown in Fig. 9.a with the fourth and the fifth stages shown in Figs. 9.b and 9.c. Here, the diseased patch is clearly learned apart from the background or the soil patches.



(a) HSL output image (b) 4c layer (c) 5c layer

Figure 9: (a) Sample output image, corresponding visualization of the leaf, and disease-relevant feature map for (b) 4c and (c) 5c layers of HSL color space detection model.

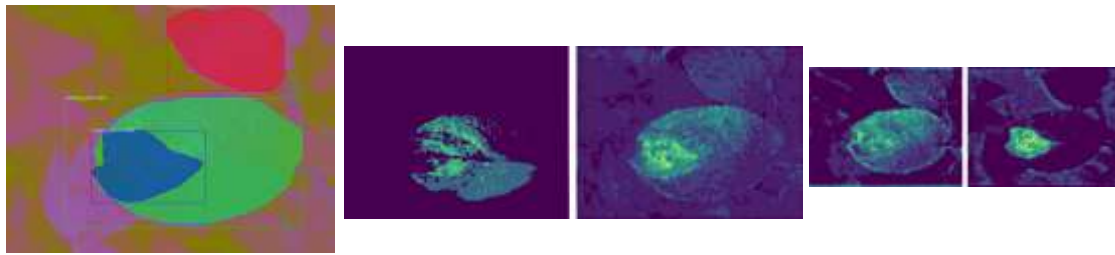
Similar to HSL, the HSV color space detection model has learned the leaf venation and disease patch structures clearly, which are visible in all the leaf feature maps of different stages for the model shown in Figs. 10.b and 10.c. Here a total of 54 and 43 feature maps have a strong correlation with the diseased patch and leaf feature respectively.



(a) HSV output image (b) 2c layer (c) 4c layer

Figure 10: (a) Sample output image, corresponding visualization of the leaf, and disease-relevant feature map for (b) 2c and (c) 4c layers of HSV color space detection model.

The sample LAB color space image shows that the disease patch (dark-blue color mask), the infected leaf (green color mask) and the healthy leaf (red color mask) have all been detected. From the feature maps shown in Fig. 11.a, it is observed that the diseased patch, the infected leaf, and the healthy leaf features are learned separately. The YCrCb color space model detects the various classes similar to the LAB color space model with a sharp distinction between the three classes.



(a) LAB output image

(b) 2c layer

(c) 3c layer

Figure 11: (a) Sample output image, corresponding visualization of the leaf, and disease-relevant feature map for (b) 2c and (c) 3c layers of LAB color space detection model.

#### 4. Discussion

The primary aim of this work is to deliver blight advisories to potato farmers in a timely and automated manner. As blight spreads fairly rapidly, farmers are advised to spray fungicide as soon as blight occurrence is detected. Therefore, in addition to successfully detecting true occurrences of blight, the blight model should also minimize the number of false positive detections. Otherwise it can lead to unnecessary spraying of fungicide and higher input costs to farmers. False positives can be a problem when using simple binary classification since the model might misinterpret the background soil patches as occurrences of blight and give a false alarm. Therefore, both a 2-class and a 4-class detection model are explored in this work.

The use of various color space transformations for preprocessing the data enables higher detection accuracy by circumventing the variations in lighting conditions on the field. While this work has presented the performance of individual color models, one may easily create a consensus system using multiple color models in parallel, to further enhance the detection. This kind of software and algorithmic approach to processing RGB images can be far more cost-effective than the use of multispectral or hyperspectral cameras. Also RGB-based data acquisition and analyses are transferable to smart-phones that are usually affordable to farmers.

The underpinnings of any successful detection model are the quality and quantity of training data. Image data in particular can vary greatly in field environments due to occlusions of the disease regions due to neighboring leaves or stems; out of focus target regions due to movement of the sensor or the target leaves themselves; illumination variation due to the season, time of the day and angle of imaging; morphological variations of leaves in terms of size, shape and texture. Therefore, a significant contribution of this work is in the collection of diverse field images of potato across different geographies and time instants to ensure a heterogeneous training data.

Modern cell phone cameras have improved in their imaging capability along with the software enhanced image processing offered by such phones. The acquisition of data using a variety of cell phones might lead to a model that can find wide applicability when many farmers are hesitant to adopt/ invest in aerial or ground-based phenotyping equipment. Apart from model performance, inference time and memory space utilization are also important metrics for smart-phone application. Therefore, the challenge will be to reduce the model size, while retaining its performance. Inference from a single image presently takes about 1-2 seconds, which makes it of practical value.

In practice, the farmer will need alerts of even a single occurrence of blight to contain it in the initial stage. In this context, it may be noted that the mAP and mAR scores provided in Table. 3 are quite conservative due to the underlying concept of IoU. While evaluating the performance of the model, a detection is considered correct only if the model is able to place a bounding box around the disease that has at least 50% area of intersection with a ground truth box demarcated by an expert. While this provides good standardization for model evaluation across various application domains, it may be noted that in the context of disease detection, the performance of the model is gauged depending on how the expert annotates the ground truth. Table. 4, on the other hand, gives a more liberal interpretation of the model performance by testing how well the model can demarcate the disease without reference to the specific ground truth annotations. Thus, the results presented in Table. 4 show that the model presented in this work can lead to more optimistic outcomes for the potato farmer.

## 5. Conclusions

This work has demonstrated a potato blight detection model using the deep learning approach that can be applied in field conditions, for aiding the farmer in making real-time decisions. In order to improve the detection performance of the model on data acquired from easily available RGB sensors, the input data are mathematically transformed to other color spaces to aid the training of the Mask R-CNN model. It is observed that training in the LAB color space provides the highest performance metrics with 80.1% and 81.4% mAP for the 2-class and 4-class detection models, respectively. The XYZ color space has the lowest mAP values for both detection models, yielding 77.4% and 70.8%, respectively. However, the model can provide an optimistic performance of ~98% overall precision for disease detection in the real-world scenario. The feature maps of intermediate layers of the trained detection models were observed and it was found that color spaces with better performance enabled the model to learn fine features of the disease patch, the leaf and the soil patch such as color, texture, leaf venation and leaf shape. The inference time per image and size of the detection models allow quick response when deployed in the field. This work could be extended to gauge the disease severity by quantifying the number and the size of blight disease patches per leaf.

## Acknowledgments:

**General:** The authors express their gratitude to Dr. Sanjay Rawal, Dr. Prince Kumar, Portia D Singh, Krishan Kumar, Mahesh Vikal, Kawalpreet and Harish Kumar for data collection and management.

**Author contributions:** JJ carried out the simulations and generated the results with assistance from GS. GS curated and annotated the datasets. SS conceived, designed and supervised the work. JJ and SS wrote the manuscript. S.Sharma and JS supervised the annotation and classification of diseases. SKM and VKD supervised the data management and testing.

**Funding:** This research was supported by the Government of India's Department of Biotechnology under the FarmerZone<sup>TM</sup> initiative (# BT/IN/Data Reuse/2017-18) and the Ramalingaswami Re-entry fellowship (# BT/RLF/Re-entry/44/2016).

**Competing interests:** SS is an advisor to Arnetta Technologies Pvt. Ltd, a startup venturing into breeding management systems for cereal crops. The authors declare that there is no conflict of interest regarding the publication of this article.

**Data Availability:** All data used to train and test the model presented in this paper is freely available upon request.

558 **References**

- 559 [1] Arora, R. K., Sanjeev Sharma, and B. P. Singh, "Late blight disease of potato and its  
560 management," *Potato Journal*, vol. 41, no. 1, pp. 16–40, 2014.
- 561 [2] Haverkort, A. J., P. C. Struik, R. G. F. Visser, and E. J. P. R. Jacobsen, "Applied  
562 biotechnology to combat late blight in potato caused by *Phytophthora infestans*," *Potato  
563 Research*, vol. 52, no. 3, pp. 249-264 , 2009.
- 564 [3] M. Islam, A. Dinh, K. Wahid, and P. Bhowmik, "Detection of potato diseases using image  
565 segmentation and multiclass support vector machine," In 30th Canadian conference on  
566 electrical and computer engineering (CCECE), pp. 1-4. IEEE, 2017.
- 567 [4] Vibhute, Anup, and Shrikant K. Bodhe, "Applications of image processing in agriculture: a  
568 survey," *International Journal of Computer Applications*, vol. 52, no. 2, 2012.
- 569 [5] Barbedo, Jayme Garcia Arnal, "Plant disease identification from individual lesions and spots  
570 using deep learning," *Biosystems Engineering*, vol. 180, pp. 96-107, 2019.
- 571 [6] S. Parkes, S. Teltscher, ITU, UI. "ICT Facts and Figures–The world in 2015," 2015.
- 572 [7] Mohanty, Sharada P., D. P. Hughes, and M. Salathé, "Using deep learning for image-based  
573 plant disease detection," *Frontiers in plant science*, vol. 7, p. 1419, 2016.
- 574 [8] Biswas, Sandika, Bhushan Jagyasi, Bir Pal Singh, and Mehi Lal, "Severity identification of  
575 Potato Late Blight disease from crop images captured under uncontrolled environment," *IEEE  
576 International Humanitarian Technology Conference(IHTC)*, IEEE, pp. 1-5, 2014.
- 577 [9] Krizhevsky, Alex, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep  
578 convolutional neural networks," *Communications of the ACM*, vol. 60, no. 6, pp. 84-90, 2017.
- 579 [10] Y. Toda, F. Okura, "How Convolutional Neural Networks Diagnose Plant Disease", *Plant  
580 Phenomics*, vol. 2019, Article ID: 9237136, 2019.
- 581 [11] LeCun, Yann, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to  
582 document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278-2324, 1998.
- 583 [12] Simonyan, Karen, and A. Zisserman, "Very deep convolutional networks for large-scale  
584 image recognition," *arXiv preprint, arXiv:1409.1556*, 2014.
- 585 [13] Szegedy, Christian, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V.  
586 Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," In *Proceedings of the IEEE  
587 conference on computer vision and pattern recognition*, pp. 1-9, 2015.
- 588 [14] He. K, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," In  
589 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770-778,  
590 2016.
- 591 [15] O'Mahony, N., Campbell, S., Carvalho, A., Harapanahalli, S., Hernandez, G.V.,  
592 Krpalkova, L., Riordan, D. and Walsh, J., "Deep learning vs. traditional computer vision," In  
593 *Science and Information Conference . Springer, Cham*, pp. 128-144, 2019.
- 594 [16] Ferentinos, Konstantinos P, "Deep learning models for plant disease detection and  
595 diagnosis," *Computers and Electronics in Agriculture*, vol. 145, pp. 311-318, 2018.
- 596 [17] Arsenovic, M. Karanovic, S. Sladojevic, A. Anderla, and D. Stefanovic, "Solving current  
597 limitations of deep learning based approaches for plant disease detection," *Symmetry*, vol. 11,  
598 no. 7, p. 939, 2019.
- 599 [18] Girshick. R, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate  
600 object detection and semantic segmentation," In *Proceedings of the IEEE conference on  
601 computer vision and pattern recognition*, pp. 580-587, 2014.
- 602 [19] R. Joseph, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-  
603 time object detection," In *Proceedings of the IEEE conference on computer vision and pattern  
604 recognition*, pp. 779-788, 2016.



- 605 [20] L. Wei, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C-Y Fu, and A. C. Berg, "Ssd: Single  
606 shot multibox detector," In European conference on computer vision, pp. 21-37, Springer,  
607 Cham, 2016.
- 608 [21] S. Ren, K. He, R. Girshick, J. Sun, "Faster R-CNN: Towards Real-Time Object Detection  
609 with Region Proposal Networks," IEEE Transactions on Pattern Analysis and Machine  
610 Intelligence, vol. 39, pp. 1137–1149, 2017.
- 611 [22] K. He, G. Gkioxari, P. Dollar, R. Girshick, "Mask R-CNN," In 2017 IEEE International  
612 Conference on Computer Vision (ICCV), IEEE, pp. 2980–2988, 2017.
- 613 [23] S. Zhang, C. Zhang, X. Wang, Y. Shi, "Cucumber leaf disease identification with global  
614 pooling dilated convolutional neural network," Computers and Electronics in Agriculture, vol.  
615 162, pp. 422–430, 2019.
- 616 [24] L. A. da Silva, P. O. Bressan, D. N. Goncalves, D. M. Freitas, B. B. Machado, W. N.  
617 Goncalves, "Estimating soybean leaf defoliation using convolutional neural networks and  
618 synthetic images," Computers and Electronics in Agriculture, vol. 156, pp. 360–368, 2019.
- 619 [25] M. D. Zeiler, R. Fergus, "Visualizing and Understanding Convolutional Networks," In  
620 European Conference on Computer Vision, Springer, Cham, pp. 818–833, 2014.
- 621 [26] S. H. Lee, C. S. Chan, S. J. Mayo, P. Remagnino, "How deep learning extracts and learns  
622 leaf features for the plant classification," Pattern Recognition, vol. 71, pp. 1–13, 2017.
- 623 [27] N. A. Ibraheem, M. M. Hasan, R. Z. Khan, P. K. Mishra, "Understanding Color Models:  
624 A Review," ARPN Journal of Science and Technology, vol. 2, 2012.
- 625 [28] H. K. Kim, J. H. Park, H. Y. Jung, "An Efficient Color Space for Deep-Learning Based  
626 Traffic Light Recognition," Journal of Advanced Transportation, pp. 1–12, 2018.
- 627 [29] D. Khattab, H. M. Ebied, A. S. Hussein, M. F. Tolba, "Color image segmentation based  
628 on different color space models using automatic GrabCut.," The Scientific World Journal, vol.  
629 126025, 2014.
- 630 [30] Das, S., Roy, D. and Das, P., "Disease Feature Extraction and Disease Detection from  
631 Paddy Crops Using Image Processing and Deep Learning Technique," In Computational  
632 Intelligence in Pattern Recognition, Springer, pp. 443-449, 2020.
- 633 [31] Ma, J., Du, K., Zhang, L., Zheng, F., Chu, J. and Sun, Z., "A segmentation method for  
634 greenhouse vegetable foliar disease spots images using color information and region growing,"  
635 Computers and Electronics in Agriculture, vol. 142, pp. 110-117, 2017.
- 636 [32] T. Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, C. L. Zitnick,  
637 "Microsoft coco: Common objects in context," In European conference on computer vision,  
638 Springer, pp. 740–755, 2014.
- 639 [33] Farmerzone-website, URL: <http://www.farmerzone.in/>, 2018.
- 640 [34] K. N. Plataniotis, A. N. Venetsanopoulos, Color image processing and applications.  
641 Springer Science & Business Media, 2013.
- 642 [35] OpenCV: Color conversions, URL: [https://docs.opencv.org/3.4.0/de/d25/imgproc\\_color\\_conversions.html](https://docs.opencv.org/3.4.0/de/d25/imgproc_color_conversions.html), 2017.
- 643 [36] Lampert, T. A., A. Stumpf, and P. Gançarski., "An empirical study into annotator  
644 agreement, ground truth estimation, and algorithm evaluation," IEEE Transactions on Image  
645 Processing, vol. 25, no. 6, pp. 2557-2572, 2016.
- 646 [37] McHugh, Mary L., "Interrater reliability: the kappa statistic," Biochemia medica:  
647 Biochemia medica, vol. 22, no. 3, pp. 276-282, 2012.
- 648 [38] Y. Kim, FasterRCNN, URL:<https://github.com/you359/Keras-FasterRCNN>, 2017.
- 649 [39] F. Islam, M. N. Hoq and C. M. Rahman, "Application of Transfer Learning to Detect Potato  
650 Disease from Leaf Image," In IEEE International Conference on Robotics, Automation,  
651 Artificial-intelligence and IoT (RAAICON), pp. 127-130, 2019.
- 652

- 653 [40] D. Tiwari, M. Ashish, N. Gangwar, A. Sharma, S. Patel and S. Bhardwaj, "Potato Leaf  
654 Diseases Detection Using Deep Learning," In International Conference on Intelligent  
655 Computing and Control Systems (ICICCS), pp. 461-466, 2020.
- 656 [41] A. Dutta, A. Zisserman, "The VIA Annotation Software for Images," Audio and Video,  
657 arXiv preprint arXiv:1904.10699, 2019.
- 658 [42] K. M. Ting, "Confusion matrix, Encyclopedia of Machine Learning and Data Mining," pp.  
659 260–260, 2017.
- 660 [43] M. Everingham, L. V. Gool, C. K. Williams, J. Winn, A. Zisserman, "The pascal visual  
661 object classes (voc) challenge," International journal of computer vision, vol. 88, pp. 303–338,  
662 2010.
- 663 [44] COCO-website, URL: <http://cocodataset.org>, 2020.
- 664 [45] E. Hadjidemetriou, M. Grossberg, S. Nayar, "Multiresolution histograms and their use for  
665 recognition," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 26, pp.  
666 831–847, 2004.
- 667 [46] OpenCV: Histograms-1:Find,Plot,Analyze., URL: [https://docs.opencv.org/3.1.0/d1/db7/  
668 tutorial\\_py\\_histogram\\_begins.html](https://docs.opencv.org/3.1.0/d1/db7/tutorial_py_histogram_begins.html), 2015.