

# DeepCount: In-Field Automatic Quantification of Wheat Spikes Using Simple Linear Iterative Clustering and Deep Convolutional Neural Networks

Pouria Sadeghi-Tehran<sup>1\*</sup>, Nicolas Virlet<sup>1</sup>, Eva M. Ampe<sup>2</sup>, Piet Reyns<sup>2</sup>, Malcolm J. Hawkesford<sup>1</sup>

<sup>1</sup>Rothamsted Research (BBSRC), United Kingdom, <sup>2</sup>Limagrain (France), France

*Submitted to Journal:*  
Frontiers in Plant Science

*Specialty Section:*  
Technical Advances in Plant Science

*Article type:*  
Original Research Article

*Manuscript ID:*  
466841

*Received on:*  
18 Apr 2019

*Revised on:*  
11 Jul 2019

*Frontiers website link:*  
[www.frontiersin.org](http://www.frontiersin.org)

### *Conflict of interest statement*

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest

### *Author contribution statement*

P.S.T proposed and developed the computer vision methods. P.S.T conducted the image processing analysis. N.V performed the statistical analysis. N.V planned and conducted the field experiments under the Scanalyzer. M.J.H contributed to the revision of the manuscript and supervised the project. All authors gave final approval for publication.

### *Keywords*

wheat ear counting, crop yield, deep learning in agriculture, Semantic segmentation, Superpixels, phenotyping, automated phenotyping system, Machine learning in agriculture

### *Abstract*

Word count: 314

Crop yield is an essential measure for breeders, researchers and farmers and is comprised of and may be calculated by the number of ears/m<sup>2</sup>, grains per ear and thousand grain weight. Manual wheat ear counting, required in breeding programmes to evaluate crop yield potential, is labour intensive and expensive; thus, the development of a real-time wheat head counting system would be a significant advancement.

In this paper, we propose a computationally efficient system called DeepCount to automatically identify and count the number of wheat spikes in digital images taken under the natural fields conditions. The proposed method tackles wheat spike quantification by segmenting an image into superpixels using Simple Linear Iterative Clustering (SLIC), deriving canopy relevant features, and then constructing a rational feature model fed into the deep Convolutional Neural Network (CNN) classification for semantic segmentation of wheat spikes. As the method is based on a deep learning model, it replaces hand-engineered features required for traditional machine learning methods with more efficient algorithms.

The method is tested on digital images taken directly in the field at different stages of ear emergence/maturity (using visually different wheat varieties), with different canopy complexities (achieved through varying nitrogen inputs), and different heights above the canopy under varying environmental conditions. In addition, the proposed technique is compared with a wheat ear counting method based on a previously developed edge detection technique and morphological analysis. The proposed approach is validated with image-based ear counting and ground-based measurements. The results demonstrate that the DeepCount technique has a high level of robustness regardless of variables such as growth stage and weather conditions, hence demonstrating the feasibility of the approach in real scenarios.

The system is a leap towards a portable and smartphone assisted wheat ear counting systems, results in reducing the labour involved and is suitable for high-throughput analysis. It may also be adapted to work on RGB images acquired from UAVs.

### *Contribution to the field*

Dear Dr. Andreas Hund, We would like to submit our original article entitled, "DeepCount: In-Field Automatic Quantification of Wheat Spikes Using Simple Linear Iterative Clustering and Deep Convolutional Neural Network", for consideration for publication in Frontiers. The authors of this manuscript are Pouria Sadeghi-Tehran, Nicolas Virlet, Eva Ampe, Piet Reyns, and Malcolm J. Hawkesford. We declare that the manuscript has not been submitted for publication elsewhere. This manuscript presents an automated model for identifying wheat spikes under natural field conditions based on a completely data-driven framework. We utilised computer vision approaches known as simple linear iterative clustering and convolutional neural networks to achieve high quality results. The proposed model can adapt to various environmental challenges faced in the field conditions and robust enough to be used as a high-throughput post-processing method to quantify the number of spikes for large-scale breeding programs. In addition, the performance of the proposed methods is compared with a state-of-the-art image processing technique in various environmental conditions. This work represents an advancement in the development of computer vision tools for application on field grown wheat canopies. Furthermore, although the proposed method is primarily focused on wheat ear counting, it could also be transferred to other applications such as identifying weeds, diseases, etc. Yours Sincerely, Pouria Sadeghi-Tehran

### *Funding statement*



Rothamsted Research receives support from the Biotechnology and Biological Sciences Research Council (BBSRC) of the UK as part of the Designing Future Wheat (BBS/E/C/00010220) and Defra Wheat Genetic Improvement Network (WGIN) (CH1090) projects.

### *Ethics statements*

(Authors are required to state the ethical considerations of their study in the manuscript, including for cases where the study was exempt from ethical approval procedures)

*Does the study presented in the manuscript involve human or animal subjects:* No

### *Data availability statement*

Generated Statement: All datasets generated for this study are included in the manuscript and the supplementary files.

In review

# DeepCount: In-Field Automatic Quantification of Wheat Spikes Using Simple Linear Iterative Clustering and Deep Convolutional Neural Networks

1 Pouria Sadeghi-Tehran<sup>1\*</sup>, Nicolas Virlet<sup>1</sup>, Eva M. Ampe<sup>2</sup>, Piet Reynolds<sup>2</sup>, Malcolm J. Hawkesford<sup>1</sup>

2 <sup>1</sup>Plant Sciences Department, Rothamsted Research, Harpenden, United Kingdom

3 <sup>2</sup>Phenotyping, Near Infrared and Research Automation Group, Limagrain Europe, Netherlands

4 **\* Correspondence:**

5 Pouria Sadeghi-Tehran

6 pouria.sadeghi-tehran@rothamsted.ac.uk

7 **Keywords: wheat ear counting, crop yield, deep learning in agriculture, semantic segmentation,**  
8 **superpixels, phenotyping, automated phenotyping system**

9 **Abstract**

10 Crop yield is an essential measure for breeders, researchers and farmers and is comprised of and may  
11 be calculated by the number of ears/m<sup>2</sup>, grains per ear and thousand grain weight. Manual wheat ear  
12 counting, required in breeding programmes to evaluate crop yield potential, is labour intensive and  
13 expensive; thus, the development of a real-time wheat head counting system would be a significant  
14 advancement.

15 In this paper, we propose a computationally efficient system called *DeepCount* to automatically  
16 identify and count the number of wheat spikes in digital images taken under the natural fields  
17 conditions. The proposed method tackles wheat spike quantification by segmenting an image into  
18 superpixels using Simple Linear Iterative Clustering (SLIC), deriving canopy relevant features, and  
19 then constructing a rational feature model fed into the deep Convolutional Neural Network (CNN)  
20 classification for semantic segmentation of wheat spikes. As the method is based on a deep learning  
21 model, it replaces hand-engineered features required for traditional machine learning methods with  
22 more efficient algorithms.

23 The method is tested on digital images taken directly in the field at different stages of ear  
24 emergence/maturity (using visually different wheat varieties), with different canopy complexities  
25 (achieved through varying nitrogen inputs), and different heights above the canopy under varying  
26 environmental conditions. In addition, the proposed technique is compared with a wheat ear counting  
27 method based on a previously developed edge detection technique and morphological analysis. The  
28 proposed approach is validated with image-based ear counting and ground-based measurements. The  
29 results demonstrate that the *DeepCount* technique has a high level of robustness regardless of variables  
30 such as growth stage and weather conditions, hence demonstrating the feasibility of the approach in  
31 real scenarios.

32 The system is a leap towards a portable and smartphone assisted wheat ear counting systems, results  
33 in reducing the labour involved and is suitable for high-throughput analysis. It may also be adapted to  
34 work on RGB images acquired from UAVs.

35 **1 Introduction**

36 Yield is composed of three components: number of ears per unit area, number of grains per ear, and  
37 grain weight, some which may be estimated during the growing season. The early estimation of pre-  
38 harvest yield allows breeders more rapid germplasm assessment and enables farmers to adjust  
39 cultivation practices to optimise production. Manual counting protocols have been the only way of  
40 calculating the number of ears per square metre (ears/m<sup>2</sup>). Breeders can identify and count wheat spikes  
41 visually, however; manual counting of wheat spikes is labour intensive and time-consuming. In  
42 addition, these tasks may need to be performed on many thousands of cultivars, which is likely to  
43 introduce human-error into the obtained data. An ideal alternative would be the development of  
44 automated systems operating under field conditions. Recent advances on automated data acquisition  
45 systems (Busemeyer et al., 2013; Kirchgessner et al., 2017; Virlet et al., 2016), allow a high spatial  
46 sampling due to the rapidity of the image acquisition process which enables all possible measurements  
47 of crop growing status. Even though the ability to acquire data is relatively fast and easy, challenges  
48 remain in terms of the data mining of images. Computer vision offers an effective choice for analysing  
49 high-throughput image-based phenotyping due to low-cost (relative to man-hours invested into manual  
50 observations) and the requirement for minimal human intervention. Although current computer vision  
51 systems are increasingly powerful and capable, they still need to overcome the difficulties associated  
52 with images acquired under field conditions. Environmental noise causes major challenges for  
53 computer vision-based techniques in identifying features-objects of interest such as wheat spikes ~~under~~  
54 ~~natural field conditions.~~ For exampleSome challenges include; (i) plant movements and/or stability of  
55 handheld cameras may cause blurred images (ii) dark shadows or sharp brightness may appear in  
56 images due to natural condition and light variations in the field even though a camera is set to auto  
57 exposure (iii) overlaps between ears due to a floppy attitude of the ears may also cause additional  
58 difficulties, especially with the presence of awns in some cultivars, and (iv) Moreover, spikes in  
59 different varieties change significantly through ~~the~~ development stages, as ~~the~~ spikes show ~~the~~ only  
60 little resemblance-similarity between the early and later growth stages.

61 Several studies have utilised image-based automatic wheat ear counting for early evaluation of yields  
62 (Cointault et al., 2008; Cointault and Gouton, 2007; Fernandez-Gallego et al., 2018). These methods,  
63 have mainly used/relies on image data extraction techniques ~~that were~~ related to characteristics of  
64 colour, texture, and morphological operations. Cointault *et. al* (2008) proposed a mobile platform to  
65 acquire data where visible images were taken by a digital camera located vertically above the field of  
66 view using a tripod. The field of view is a closed system delimited by a black matte frame to control  
67 variabilities in illumination and weather conditions. The proposed framework creates a homogeneous  
68 environment and blocks unwanted image effects. Subsequently, the authors improved their platform  
69 by collecting images in different lighting conditions without any structure blocks (Cointault et al.,  
70 2008). The main drawback is the restricted data acquisition pipeline required for the system to operate.  
71 For instance, prior knowledge of the environment is required to achieve an optimum result; moreover,  
72 even with the current restrictions only a small number of images were selected based on which the  
73 authors felt presented "good illumination". In a similar approach (Cointault et al., 2008a; Cointault and  
74 Gouton, 2007; Fernandez-Gallego et al., 2018), a supervised classification method was proposed to  
75 distinguish three classes of leaves, soil and ears. In the end, morphological operations were applied for  
76 counting the number of blobs (potentially ears) from the binary image with the pre-assumptions of the  
77 shapes of the ears. Each pixel is represented by colour and texture properties. As suggested, a hybrid  
78 space is constructed to address a sensitivity of colour properties to the intensity variations in an image.  
79 The method has been tested on a limited number of wheat varieties without awns with a low level of  
80 wheat ear density; ~~moreover~~nonetheless, no evaluation was carried out to validate the accuracy of the  
81 proposed method with the manual measurements. In another study, Fernandez *et. al.* (2018) applied a

82 Fourier filtering and two-dimensional discrete Fast Fourier transform (FFT) ([Cooley and Tukey, 1965](#))  
83 to distinguish wheat ears from the background. The approach performs, in three main steps of high-  
84 pass filtering, thresholding and mathematical morphology, operations to eliminate "non-wheat" pixel  
85 groups which are small and scattered. The threshold is pre-defined by a user to determine if pixels  
86 should be identified as foregrounds (ears) or background (leaf, soil, etc.). The drawback is that a wrong  
87 choice of the threshold value may result in distortion and low performance of the whole system in  
88 different environments. Finally, Zhou *et. al* 2018 proposed a twin-support-vector machine  
89 segmentation method to segment wheat ears from visible images. The method relies on the hand-  
90 engineered features including colour, texture, and edge histogram descriptor. The images were  
91 collected from the side at 45° above the horizontal because colour and texture were suggested being  
92 typically more substantial from this perspective.

93 At the core, the success of any of the current state-of-the-art methods crucially depends on the feature  
94 representation of the images. While the aforementioned methods use hand-crafted features to represent  
95 images by encoding of various features including corners, edges, texture and colour schemes, the  
96 features are tailored to a specific condition and their effectiveness are inherently limited as these  
97 approaches mainly operate at the primitive level. Unlike conventional feature extraction techniques,  
98 which often use shallow architecture and solely rely on human-crafted features, relatively new  
99 learning-based methods based on Convolutional Neural Networks (CNNs) show promising results for  
100 visual analysis. CNN models attempt to model high-level abstractions in images by employing deep  
101 architectures composed of multiple non-linear transformations ([Lomonaco, 2015; Schmidhuber, 2015](#)).  
102 In CNN, features are extracted at multiple levels and allow the system to learn complex functions that  
103 directly map raw sensory input data to the output, without relying on hand-engineered features using  
104 domain knowledge. The convolution is an operation of applying the filter on a single colour image to  
105 enhance some of its features. One-to-one convolutions take a single image as an input and return a  
106 single image as an output. However, in CNN different kinds of convolutions exist. For instance, in one-  
107 to-many convolutions, a single input image is passed to  $k$  filters; then each filter is used to generate a  
108 new output image. Alternatively, in many-to-many convolutions, there are  $n$  inputs and  $m$  outputs  
109 where each output image is connected to one or more input image characterised by  $k$  filters ([Lomonaco,  
110 2015](#)). Potentially, this capability makes the deep neural network more robust to different types of  
111 variations in digital images. As a result, the model can adapt to such differences and has the capacity  
112 to learn complex models.

113 In recent years, CNNs have shown usefulness in a large variety of natural language processing and  
114 computer vision applications, including segmentation and image classification, and often surpassed  
115 the-state-of-the-art techniques ([Krizhevsky et al., 2012; Lomonaco, 2015; Mikolov et al., 2013](#)).  
116 Despite the promising outcomes of deep learning in computer vision, there are some limitations in  
117 implementing a deep neural network. Deep learning approaches are usually computationally intensive,  
118 and their performance relies on the quantity and quality of training datasets. In most cases, in order for  
119 deep learning to show great advantages, training datasets of tens of thousands to millions are required  
120 ([Deng et al., n.d.; Ubbens et al., 2018](#)). Having a large training dataset provides deep learning models  
121 with extensive variety, which leads to an effective learned representation as a result. [Deep Neural  
122 Networks \(DNN\)](#) is an area of active research and applications to plant research are still in the early  
123 stages. There are few deep learning applications successfully applied in the field of image-based plant  
124 phenotyping ([Madec et al., 2019; Pound et al., 2017](#)). The small body of existing applications includes  
125 plant disease detection on leaf images ([Mohanty et al., 2016](#)), rice panicle segmentation ([Xiong et al.,  
126 2017](#)), leaf counting in rosette plants ([Ubbens et al., 2018](#)), wheat ear counting ([Madec et al., 2019](#)),  
127 and localising root and shoot tips ([Pound et al., 2017](#)).

128 This study utilises a novel visual-based approach based on linear iterative clustering and deep  
129 convolutional neural networks to identify and count the number of wheat spikes. The proposed method  
130 can also calculate the number of wheat ears/m<sup>2</sup> when a ground standard is present within the image.  
131 The proposed method, called *DeepCount*, alleviates the limitations and lack of separability inherent in  
132 existing wheat ear counting methods and minimise the constraints of capturing digital images taken  
133 under natural outdoor environments. The approach presented will pave the way for computationally  
134 efficient and significantly faster approaches compared to the manual techniques, leading to reducing  
135 the labour involved and enabling high-throughput analysis.

## 136 2 Materials and Methodology

137 In this study, we explore the feasibility of automatically identifying wheat spikes under natural in field  
138 conditions based on a completely data-driven framework. The main contributions of the work can be  
139 summarised as follows:

- 140 • Building high-quality dataset of annotated spikes and utilising them to train our convolutional  
141 neural network model
- 142 • Developing a deep learning model called *DeepCount* that can learn from the training dataset  
143 and then identify and segment spikes from different wheat cultivars (awns and no awns).
- 144 • Demonstrating that the constructed model can automatically quantify the number of spikes ~~in~~  
145 within visible images ~~in-under~~ natural field environments; ~~also~~, calculate the number of ears/m<sup>2</sup>  
146 when a ground standard is present.

147 Quantification of spikes may be achieved in two ways. One approach is localisation/detection of spikes,  
148 which provides not only the prediction for the whole image but also additional information regarding  
149 the spatial location of the spikes. Another technique is semantic segmentation (pixel-wise  
150 segmentation) which understands an image at pixel level. It enables dense predictions inferring labels  
151 of every pixel in the image, so that each pixel is labelled as an ear or background. Inspired by the  
152 success of the recent deep learning algorithms in computer vision applications, we propose a CNN  
153 approach combined with a superpixels technique known as simple linear iterative clustering (SLIC)  
154 ([Achanta et al., 2010](#)). The core idea is to overcome the computational complexity by using SLIC to  
155 generate homogeneous regions instead of processing at a pixel level. The homogeneous regions  
156 generated by SLIC will contain more information about the colour and texture and are less sensitive to  
157 noise as opposed to pixel-level analysis. It also reduces the complexity of subsequent ear detection and  
158 localisation tasks. The generated regions are later used as input data for the convolutional neural  
159 networks. The network is not only capable to recognise spikes but also delineate the boundaries of each  
160 spike with the canopy based on dense pixel level predictions. Figure 1 illustrates an end-to-end  
161 wheatear quantification including the offline training and online ear segmentation and counting. In the  
162 following section, we will describe the data collection/annotation process and the model architecture  
163 developed to localise wheat spikes within images and quantify them.

### 164 2.1 Experimental materials

165 The experiments were carried out at Rothamsted Research, UK (51°48'34.56"N, 0°21'22.68"W) in two  
166 fields, Great Field (Field Scanalyzer area) and Black Horse. Two experiments were conducted ~~underon~~  
167 the Field Scanalyzer platform ([Virlet et al., 2016](#)) during the growing season in 2014-2015 (hereafter  
168 referred to as 2015-FS data set) and 2015-2016 (hereafter referred to as 2016-FS data set). Six wheat  
169 cultivars (*Triticum aestivum* L. cv. Avalon, Cadenza, Crusoe, Gatsby, Soissons and Maris Widgeon)  
170 were sown on 6<sup>th</sup> November 2014 and 20<sup>th</sup> October 2015 at a planting density of 350 seeds/m<sup>2</sup>. Nitrogen



171 (N) treatments were applied as ammonium nitrate in the spring, at rates of 0 kgN.ha<sup>-1</sup> (residual soil N;  
172 N1) 100 kgN.ha<sup>-1</sup> (N2) and 200 kgN.ha<sup>-1</sup> (N3) for both years and 350 kgN.ha<sup>-1</sup> (N4, 2015-FS only).  
173 The plot sizes were 3m × 1m in 2015-FS and 2m × 1m in 2016-FS.

174 The third experiment has been funded by DEFRA since 2008, known as WGIN (Wheat Genetic  
175 Improvement Network), to provide genetic and molecular resources for research in other DEFRA  
176 projects and for a wide range of wheat research projects in the UK. In this study, we collected images  
177 from the 2015-2016 experiment (hereafter referred to as 2016-WGIN data set) at Black Horse field. 30  
178 wheat cultivars were grown at four nitrogen fertiliser treatments (N1, N2, N3 and N4), sown on 12<sup>th</sup>  
179 October 2015. Each repetition consists of a 9m × 3m “main plot”, and a 2.5m × 3m “sampling plot”,  
180 used for non-destructive measurement and destructive sampling respectively. The three experiments in  
181 this study use a split plot design (with three blocks) and were managed by local agronomic practices.

## 182 **2.2 Image acquisition**

183 The images were ~~taken-acquired~~ under conditions of natural illumination ~~conditions~~-at ~~different~~  
184 multiple stages of ear maturation with different canopy complexities achieved through ~~different-varied~~  
185 nitrogen inputs. The tests were carried out in extreme lightning conditions with typical environmental  
186 challenges faced in the field for images taken by different cameras and optics with no direct scaling  
187 relationships. Table 1 summaries the characteristics of three trials carried out in this study. The camera  
188 models include different types of commercially available visible cameras with various spatial  
189 resolutions and configurations (Table 1).

190 The images for 2015-FS and 2016-FS were collected ~~using-by~~ the Scanalyzer onboard visible camera  
191 (colour 12-bit Prosilica GT3300) at a resolution of 3,296×2,472 pixels. The camera is positioned  
192 perpendicular to the ground and was set up at a fixed distance to the ground (3.5m) for the 2015-FS  
193 experiment and at a fixed distance to the top of the canopy (2.5m) for the 2016-FS. The camera is set  
194 up in auto-exposure mode, to compensate for outdoor lighting changes.

195 In the 2016-WGIN experiment, two hand-held cameras, Canon G12 and Sony Nex-7, were used to  
196 acquire visible images with the resolution of 3,648×2,736, and 6,000×3,376 pixels, respectively (Table  
197 1). Similarly, to the Field Scanalyzer, the cameras were set up in an auto-exposure mode and held  
198 vertically over the canopy. In addition, a rapid and easy ground standard system was implemented by  
199 placing an A4 sheet over the canopy in the field of view of the camera lens (Figure 2.B). The ground  
200 system was used to transform the total number of wheat ears ~~in-within~~ an image into the number of  
201 ears/m<sup>2</sup>.

## 202 **2.3 Ground-truthingEvaluation**

203 Two different ~~ground-truthingevaluation~~ methods were ~~implemented-used~~ and compared with the  
204 automatic ear counting techniques. The first method is based on manual image-based annotation in  
205 which ears are manually ~~marked-counted~~ on the images acquired by the Field Scanalyzer platform  
206 (2015-FS and 2016-FS data sets). Wheat ears were interactively marked using the VIA image annotator  
207 (Dutta et al., n.d.), which enabled the automatic printing of the incremental number on each individual  
208 ear.

209 The second ground-truthing method is based on field manual measurements carried out for all three  
210 experiments. In the 2015-FS and 2016-FS experiments, ears were manually counted on six rows of 1  
211 m length, corresponding to 1 m<sup>2</sup> area, for each plot. In the 2016-WGIN trial, the number of ears/m<sup>2</sup>  
212 were estimated based on the method presented in Pask et. al (2012a). Samples of 4 rows of 1 metre

length ~~was done~~ ~~carried out~~ ~~were cut~~ at anthesis, ~~then and~~ the ears/m<sup>2</sup> were derived from the above-ground biomass (ABG) and the dry weight (DW) of the fertile culm:

$$\text{Ears/m}^2 = \text{AGB (g/m}^2) / \text{DW\_fertile culm (g)}$$

Figure 2 shows the representation of digital images of different wheat traits taken under the Field Scanalyzer platform (Figure 2.C) and a handheld DSLR camera (Figure 2 A&B). As depicted in the sample images, the data was collected in different weather conditions, with illumination changes, ~~and~~ ~~from~~ cultivars with differences in ear shapes and sizes.

## 2.4 Annotation and generating training dataset

The fundamental part of any supervised decision-making system such as CNN is how to specify the output based on a given set of inputs or training dataset. In practice, hundreds or even thousands of annotated training datasets are required to make a good training of CNN. Even though high-throughput image-based plant phenotyping systems like Field Scanalyzer (Virlet et al., 2016) exist and generate a huge amount of image data daily, a large set of annotated images with ground-truth are not widely accessible yet within the plant phenotyping community.

To expose our CNN model to a wider variety of images, the data were collected by a hand-held DSLR Canon Camera with a resolution of 5760×3840 pixels from diverse Limagrain field trials at different stages from heading to maturation under different ambient illumination condition. The broad range of images enabled the constitution of “strong” training data set covering the ears development from multiple wheat varieties making the detection model more robust and thereby increasing the precision of the wheat spikes quantification. The graphical image annotation tool, VGG image annotator (VIA) (Dutta et al., n.d.), was used to draw boxes around the background, such as leaf, soil and soil (Figure 3.C) and draw strokes using the polygon tool around ears (Figure 3 A&B). Here, 330 representative wheat images are selected to build the annotated training dataset, in which the illumination variations, weather conditions, wheat ears shapes, and reproductive stages are all considered. As a result, 24,938 ears and 30,639 backgrounds are manually annotated.

The next step is to combat the high expense of creating a training source with their corresponding labels. The augmentation model is constructed to simulate the illumination change by adjusting the HSV colour space and applying various transformations such as random rotation, cropping, flipping, zooming, scaling, and brightness to the images that are already in the training dataset (Figure 4). In addition, a non-linear operation known as Gamma correction (also referred to as gamma encoding or gamma compression) (Rahman et al., 2016) was applied to encode and decode luminance in the images. The augmented images are appended to the existing training samples, from which 20% of the sample set is randomly selected as the validation set (145,000 patches), and the remaining 80% is selected as the training set (580,000 patches; 300,000 ears and 280,000 backgrounds).

## 2.5 Superpixels segmentation

Most computer vision algorithms use pixel-grid as the underlying representation of an image. However, grids of pixels do not hold a semantic meaning of an image, nor represent a natural representation of a visual scene. It would be more efficient, to work with perceptually meaningful entities obtained from a low-level grouping process. Superpixel algorithms aim to group pixels into perceptually meaningful regions based on their similarity characteristics, such as colour and texture distributions. Superpixel techniques will reduce the complexity of images from thousands to millions of pixels to only a few

254 hundred superpixels; thereby, it will diminish the influence of noise and potentially improves the  
255 computational efficiency of vision algorithms.

256 In light of the fundamental importance of superpixel algorithms in computer vision, many algorithms  
257 have been proposed in the literature ([Achanta et al., 2012, 2010](#); [Li and Chen, 2015](#); [Tu et al., 2018](#)).  
258 The superpixel segmentation algorithms can be broadly categorised as graph-based segmentation and  
259 clustering-based segmentation. In graph-based techniques, an image is considered as a planar graph,  
260 where pixel vertices and pixel affinities are computed for connected pixels ([Felzenszwalb and  
261 Huttenlocher, 2004](#); [Ren and Malik, 2003](#)). Alternatively, the clustering-based method starts with a  
262 rough initial clustering of pixels, then the clusters are refined iteratively until some convergence  
263 criterion is met to form superpixels ([Achanta et al., 2010](#); [Achanta and Ssstrunk, 2017](#); [den Bergh et  
264 al., 2015](#)).

265 In this study, we use simple linear iterative clustering (SLIC) ([Achanta et al., 2012, 2010](#)), which is  
266 fast and memory efficient for generating superpixels ([Achanta et al., 2012](#)). As opposed to other  
267 superpixels algorithms with many difficult-to-tune parameters, SLIC is simple to use in which the  
268 number of desired superpixels is its sole parameter. The spectral-spatial distance is measured between  
269 each pixel to its cluster centre and then the cluster centres are updated using  $K$ -means clustering  
270 technique. For  $N$  pre-specified superpixels, clustering pixels are represented based on their colour  
271 similarity (CIELAB colour space) and pixel proximity in the 5-D space  $C_i = [l_i, a_i, b_i, x_i, y_i]$  where  $i =$   
272  $[1, N]$ . In this study, based on our experience, the number of superpixels is set to  $N = 3000$  to avoid  
273 over segmentation and to produce roughly equally sized superpixels. We can also control the trade-off  
274 between the compactness of the superpixels and boundary adherence ([Achanta et al., 2012](#)). It means  
275 SLIC can prevent small or disconnected areas or islands within a larger region (Figure 5). The candidate  
276 regions are then used as inputs for the CNN model to perform pixel-wise segmentation. Feeding the  
277 network with image descriptors extracted from the candidate regions enables the model to learn local  
278 information such as texture and shape rather than using the pixel-grids.

## 279 **2.6 Architecture of the convolutional neural network model**

280 As previously mentioned, SLIC reduces the computational complexity by partitioning an image into  
281 homogeneous regions instead of extracting features at the pixel level (Figure 5). However, the SLIC  
282 method, like many other superpixel techniques ([Felzenszwalb and Huttenlocher, 2004](#); [Li and Chen,  
283 2015](#); [Ren and Malik, 2003](#); [Wang et al., 2017](#)), relies on handcrafted features; thus, often fails to  
284 separate objects within an image in appropriate regions (Figure 5.C & Figure 6.A.1). To address the  
285 limitation, the proposed CNN model classifies each superpixel at a pixel-level as opposed to  
286 characterising the content of the entire candidate region and predict a single label. The network takes  
287 each candidate region as input data and outputs a pixel-level segmented of the region (Figure 6.A.2 &  
288 B.2).

289 In general, semantic segmentation architecture in CNN can be broadly categorised as an encoder  
290 network followed by a decoder network. The encoder network gradually reduces the spatial dimension  
291 of the input by down-sampling and developing lower-resolution feature mappings which are learned  
292 to be highly efficient at discriminating between classes. To get the dense pixel-wise classification, the  
293 decoder network semantically projects the discriminative features learnt by the encoder onto the pixel  
294 space by up-sampling the feature representations into a full-resolution segmentation map. There are  
295 usually shortcut connections from encoder to decoder to help the decoder recover the object details  
296 better.



297 In this work, we leverage an existing model known as U-Net which was originally designed for  
298 biomedical image segmentation for identifying lung nodules in a CT scan ([Ronneberger et al., 2015](#)).  
299 The U-Net architecture consists of a contracting path to capture context and an asymmetric expanding  
300 path that enables precise localisation. The model concatenates the encoder feature maps to up-sampled  
301 feature maps from the decoder at every stage. The concatenation allows the decoder at each stage to  
302 learn back relevant features that are lost when pooled in the encoder. Normally, U-Net is trained from  
303 scratch starting with randomly initialised weights (optimisation variables). Since up-sampling in the  
304 decoder is a sparse operation we need a good prior from earlier stages to better represent the  
305 localization.

306 Since transfer learning proved to be a powerful technique for semantic segmentation models such as  
307 U-Net like architectures ([Igloukov and Shvets, 2018](#)), we used a pre-trained VGG model ([Simonyan  
308 and Zisserman, 2014](#)) without fully connected layers as its encoder mechanism followed by a decoder  
309 network as the original U-Net to further improve the performance of pixel level dense classification.  
310 The VGG family of CNN can be characterised by two components. 1) all convolutional layers in the  
311 network use  $3 \times 3$  filters. 2) multiple convolutional layer sets are stacking together before applying a  
312 pooling operation. Normally the number of consecutive convolutional layers increases the deeper the  
313 network goes ([Simonyan and Zisserman, 2014](#)). The VGG-16 used in this work, was proposed by a  
314 group of researchers in Oxford and the winner of the ImageNet competition ([Deng et al., n.d.](#)) in 2013.  
315 It uses a stack of convolution layers with small receptive fields in the first layers instead of few layers  
316 with big receptive fields.

317 By using an existing architecture in which the weights are initialised on big data sets such as ImageNet,  
318 the network can converge faster and learn more general filters. To construct the encoder, the fully  
319 connected layers were removed and replaced with a single convolutional layer of 512 channels that  
320 serves as a bottleneck part of the network to separate the encoder from the decoder. The network  
321 contains a total of four max-pooling layers. For each of the pooling layers, the spatial size of the feature  
322 map is reduced by a factor of two vertically and horizontally.

323 The decoder part of the network consists of up-sample and concatenation with an output of the  
324 corresponding part of the decoder followed by regular convolution operations (Figure 8). Since the  
325 pre-trained VGG model takes an input of  $224 \times 224$  pixels with 3 channels, the irregular superpixels  
326 need to be resized to achieve a proper input into the model. The network takes superpixels as inputs  
327 and outputs a segmented version of the inputs. Each pixel is labelled as 1 (wheat spikes) or 0  
328 (background), which generated a binary image (Figure 7). After the semantic segmentation, the median  
329 filter is applied to minimise the noise and remove the result of misclassification over the binary image.  
330 In this process, a window size of seven pixels slides over the entire image, pixel by pixel. Then, the  
331 pixel values from the window are sorted numerically and replaced with a median value of neighbouring  
332 pixels. In the end, for contour quantification, a classical image processing algorithm known as the  
333 watershed technique is used for post-processing to further segmentation of individual contour.

### 334 2.6.1 Loss Function

335 The role of loss function in our parameterised learning was investigated. The parameterised learning  
336 will allow us to take sets of input data (ears and background) and their class labels and learn a function  
337 that maps the input to the output predictions by defining a set of parameters and optimising over them.  
338 At a basic level, a loss function quantifies how good or bad a given predictor is at classifying the input  
339 data in our dataset ([Harrington, 2012; Marsland, 2009](#)).

340 The binary cross-entropy loss function is used to quantify how accurate the CNN method is at  
 341 classifying the input data in our dataset (a brief overview of the cross-entropy loss function and the  
 342 calculations is provided in the supplementary data). A visualisation of the loss function plotted over  
 343 time for our model is shown in Figure 9. A visualisation of training accuracy, training loss, validation  
 344 accuracy, and validation loss plotted over time for the model is plotted after 15 epochs<sup>1</sup>. The smaller  
 345 the loss, the better a job the model/classifier is at modelling the relationship between the input data and  
 346 output class labels. As shown in Figure 9, loss starts slightly high but then decrease rapidly and  
 347 continues to stay low when trained on our dataset. As expected the usage of the pre-trained VGG  
 348 model helps the network to converge faster, as a result, we obtained 98% accuracy after only 15 epochs.  
 349 Furthermore, the training and validation curves match each other very closely, indicating there is no  
 350 issue of overfitting with the training process.

## 351 2.7 Hand-crafted features extraction techniques for wheat ear quantification

352 A hand-crafted image-based method presented in (Jansen et al., n.d.) was compared with the proposed  
 353 *DeepCount* model. The technique is based on an edge detection technique and several morphological  
 354 image processing operations. Firstly, the image is converted from a 3-D RGB image into 2-D greyscale  
 355 representation of the image (Figure 10.B), then the edge detection based on Sobel kernel (Kaufman et  
 356 al., 1994) performs a 2-D spatial gradient measurement on the grey image to emphasise regions of high  
 357 spatial frequency that correspond to edges which returns a binary image (Figure 10.C). Edges may  
 358 correspond to boundaries of an object, boundaries of shadowing or lighting conditions and/or  
 359 boundaries of parts within an object in an image. The next steps are morphological operations including  
 360 dilation to increase the size of foreground pixels (Figure 10.D), which is useful for joining broken parts  
 361 of the image. Filling the holes (Figure 10.E), removing small objects (Figure 10.F) are the fifth and  
 362 sixth steps. The final step is erosion where pixels near the boundary of an object in the image will be  
 363 discarded. A foreground pixel in the input image will be kept only if all pixels inside the structuring  
 364 element are bigger than zero; otherwise, the pixels are set to zero (Figure 10.G). In the end, a list of all  
 365 contours is returned, and their numbers are printed out on the RGB image (Figure 10.H). The hand-  
 366 crafted method will be referred hereafter as the edge method.

## 367 3 Results and discussions

368 The performance of the proposed *DeepCount* model was evaluated against the hand-engineered edge  
 369 detection method as well as two ~~ground truthing manual evaluation~~ techniques. The first technique was  
 370 based on manual counting of ears within visible images while the second ~~ground truthing evaluation~~  
 371 ~~method~~ was the field-based measurements. In addition, the ears counting performances were quantified  
 372 based on the coefficient of determination ( $R^2$ ), the root means squared error (RMSE), the relative  
 373 RMSE (rRMSE), and the bias:

$$374 \quad \text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^n (r_i - e_i)^2} \quad (1)$$

$$375 \quad \text{rRMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^n \left( \frac{r_i - e_i}{r_i} \right)^2} \quad (2)$$

---

<sup>1</sup> Epoch is a hyperparameter which is defined before training a neural network learning model. It means the learning algorithm has seen each of the training data points  $N$  times.

$$\text{Bias} = \frac{1}{N} \sum_{i=1}^n (r_i - e_i) \quad (3)$$

376

377 where  $N$  denotes the number of images,  $r_i$  and  $e_i$  are the reference and estimated counts for image  $i$ ,  
378 respectively.

379

380 The algorithm was tested on a workstation PC running a Centos7 operating system with 10-core Intel  
381 Xeon CPU, 3.6 GHz per CPU, 64 GB of memory and Nvidia Quadro M5000 video card. The CNN  
382 framework was developed in python using OpenCV library and the Keras framework. While there is  
383 no restriction in the spatial resolution of the test images, the segmentation and quantification of wheat  
384 spikes will take approximately 90-100 seconds on a single image with the resolution of  $6,000 \times 3,376$   
385 pixels. The CUDA parallel acceleration was also used to improve the processing efficiency especially  
386 for training the model. CUDA is a parallel computing platform created by NVIDIA, and the cuDNN  
387 library was developed for deep learning with GPU acceleration. The current method also has the  
388 potential to be faster in the future by CPU multithreading utilisation. ~~It should be highlighted that there~~  
389 ~~is no restriction in the spatial resolution of the tested images.~~

### 390 3.1 *DeepCount* vs. hand-crafted edge method

391 First, the performance of the ~~two~~ automatic image-based methods (*DeepCount* and the hand-crafted  
392 technique presented in section 2.7) was compared against manual image-based counting. In the image-  
393 based ~~ground-truthing evaluation~~, 33,011 ears were manually counted/~~annotated~~ from 126 images. The  
394 2015-FS and 2016-FS trials include 72 and 54 images in which 22,284 and 10,727 ears were manually  
395 counted on the images, respectively.

396 Figure 12 A&B illustrates the linear regression between the automatic methods and ~~the~~ first ~~ground-~~  
397 ~~truthing evaluation~~ method ~~for tested on~~ the 126 images. The results showed a high correlation between  
398 the automatic methods and the manual image-based counting. The *DeepCount* model has a higher  
399 coefficient of determination and lower RMSE and rRMSE ( $R^2 = 0.94$ , RMSE = 25.1, rRMSE = 11%)  
400 than the edge detection method ( $R^2 = 0.75$ , RMSE = 45.5, rRMSE = 21%) indicated that the *DeepCount*  
401 technique was closer to the visual observation. In addition, the bias values of -13.1 and -13.2 for both  
402 methods show a slight overestimation of the number of ears compared to the visual assessment (Figure  
403 12 A&B).

404 The visual inspection of the results suggested that the edge method had more false positives than the  
405 *DeepCount* model. ~~It was observed that, in some cases, where leaves or objects have clearer contrast~~  
406 ~~than their surroundings, they were misidentified as ears. In some cases, leaves or objects with clear~~  
407 ~~contrast than their surrounding were wrongly identified as ears.~~ This was expected since the edge  
408 detection is defined as discontinuities in pixel intensity, in other words, a sharp difference and change  
409 in pixel values; thus, ~~the edge detection method~~ is more prone to noise. This may also pose more  
410 difficulties for the edge method to identify ears with awns (e.g. Soissons cv). The *DeepCount* model,  
411 on the other hand, had less false positive, regardless of the cultivars or level of nitrogen. Further, visual  
412 inspection showed that the fraction of false negatives, in both automatic methods, appeared to be the  
413 failure of the watershed method to separate ears exposed to a severe degree of overlap.

414 While Fernandez-Gallego *et al* (2018) argued that the Edge method is unlikely to be reliable due to  
415 loss of RGB information during its colour transformation to grayscale, our results indicated otherwise.  
416 The edge method showed similar performances compared to the method presented by the authors. The

417 success rate metric ( $\mu$ ) used by the authors to evaluate the performance of their method showed 31.96  
 418 to 92.39% on RGB images and 65.36 to 93.01% on greyscale images, whereas we achieved a similar  
 419 range of values with 86% and 81% in the 2015-FS and 2016-FS experiments, respectively. Moreover,  
 420 the  $R^2$  values between the edge method and ~~the two ground-truthingevaluation~~ techniques (image-  
 421 based counting and ground-based measurements) are high with  $R^2 = 0.75$  and  $0.60$ , respectively (Figure  
 422 12 A&C). Nevertheless, the *DeepCount* model outperformed the edge method in every experiment  
 423 carried out in this study. Our results are also in agreement with the method presented by Madec *et al*  
 424 (2019). The authors obtained  $R^2 = 0.91$  and  $rRMSE = 5.3\%$  from their manual image-based ear  
 425 counting which is also very similar to the 2016-FS data set where the results showed  $R^2 = 0.97$  and  
 426  $rRMSE = 7\%$  (Figure S1). We also found similar outcomes between our methods and the technique  
 427 presented by Zhou *et al* (2019); however, as the performance metrics differs a quantitative comparison  
 428 is not possible.

429 Furthermore, ~~t~~The performances of the Edge and *DeepCount* methods were ~~compared to~~validated  
 430 against the ground-based measurements after ~~converting~~the numbers of ears/image were converted  
 431 into ears/m<sup>2</sup>. As shown in Figure 12 C&D, the performance degraded slightly compared to the manual  
 432 image-based ~~ground-truthingmeasurements~~ (Figure 12 A&B). ~~R<sup>2</sup> in~~ the edge method, R<sup>2</sup> reduced from  
 433  $0.75$  to  $0.60$ , whereas the performance in the *DeepCount* model dropped ~~only~~ from  $R^2 = 0.94$  to  $0.86$ .  
 434 The edge and *DeepCount* methods had a similar bias (36 and 35.3, respectively), which indicated that  
 435 both methods underestimated the number of ears/m<sup>2</sup> compared to the field data. In addition, the RMSE  
 436 increased from 45.5 to 104.9 ears/m<sup>2</sup> and 25.1 to 71.4 in both approaches, respectively.

437 A similar decrease in performance also observed in (Madec *et al.*, 2019). This is partly attributed to the  
 438 relatively different observation area used for the ground measurements and the visible images. The  
 439 spatial representativeness was therefore limited to get an accurate comparison between the automatic  
 440 counting and field-based measurements that were not measured at the same place over plots. For  
 441 instance, in the 2015-FS trial, the ground-based measurements ~~were~~ obtained from six rows  
 442 including the edge rows; however, the same area was not taken by the Field Scanalyzer. The number  
 443 of rows captured in the images varies between 3.5 to 5 rows (Figure 2.C). An additional factor may  
 444 also due to the fact that some ears are hidden deep down inside canopies or partially visible on the  
 445 borders of images which pose more difficulties for the automatic models to identify them. Further  
 446 improvement can be achieved between the automatic counting and direct counting in the field if the  
 447 same protocol is followed by both methods during data acquisition. For example, in the 2016-FS trial,  
 448 the results showed an improvement in performance ~~in the 2016-FS trial~~ when images were precisely  
 449 consistently taken from four middle rows in every plot (Table 2).

### 450 3.2 *DeepCount* model vs. field-based measurements

451 The performance of the *DeepCount* model was further evaluated against the ground-based  
 452 measurements in each individual trial and all together. As shown in Figure 13, the coefficient of  
 453 determination was higher in the 2016-FS experiment ( $R^2 = 0.89$ ) compared to the 2015-FS ( $R^2 = 0.70$ )  
 454 and 2016-WGIN ( $R^2 = 0.57$ ) trials. Also, the lowest bias was obtained in the 2016-FS (bias = 3.6)  
 455 followed by 2016-WGIN and 2015-FS with 37.4 and 59.14, respectively. As mentioned in the previous  
 456 section, the notable difference in bias between the 2016-FS and the other ~~two~~ trials may reside in the  
 457 fact that first, the measurements on the field and the visible images were obtained from the same area;  
 458 also, in the 2016-FS, the camera was set up at fixed distance to the top of a canopy (2.5m) regardless  
 459 of the height of the plots. As opposed to the 2015-FS trial where the camera was set up at a fixed  
 460 distance to the ground (3.5m) or in the 2016-WGIN trial, where the distance between the hand-held  
 461 cameras and top of canopies varies from one plot to another.



462 Furthermore, the lower performance in the 2016-WGIN trial may be associated with several factors.  
 463 ~~(i)First~~, improper placement of an A4 sheet used as a ground standard to transform the total number of  
 464 wheat ears in an image into the number of ears/m<sup>2</sup>. In order to have an accurate ear density estimation,  
 465 the sheet should be placed perpendicular to the handheld camera's viewing angle which was not the  
 466 case in many images taken from the WGIN-2016 trial. In addition, in some images, the ground standard  
 467 was partially obstructed by leaves and wheat ears. ~~(ii)Second~~, the perspective of the images may also  
 468 account for the slight lack of correlation between the proposed model and the field measurements.  
 469 While focal length does not change perspective per se, it does change how the ears are represented;  
 470 thus, it is important to capture the scene optimally. The ultra-wide angle focal length used to capture  
 471 images from 2016-WGIN (6 and 18 mm) provided a bigger field of coverage but caused a perspective  
 472 distortion particularly on the image borders. Last but not least, the~~(iii)~~ manual field measurements may  
 473 ~~likely to have~~ introduced human-error into obtained data.

474 Despite the above uncertainties, the *DeepCount* algorithm showed the same accuracy in every  
 475 experiment (rRMSE = 15%+/-1) regardless of the number of ears identified in the images (2015-FS:  
 476 309 to 655, 2016-FS: 183 to 634, 2016-WGIN: 238 to 821), types of cameras with different spatial  
 477 resolutions. The same accuracy also obtained when all three experiments were combined together (R<sup>2</sup>  
 478 = 0.72 and rRMSE = 15%). As shown in Table 1, two cameras (Canon and Sony) with different spatial  
 479 resolutions and lens focal lengths were used to acquire images. In the Canon camera, we observed  
 480 lower R<sup>2</sup> but higher bias compared to the Sony camera the images in 2016-WGIN were collected from  
 481 two hand-held cameras with different spatial resolutions and focal length. While the R<sup>2</sup> is lower and  
 482 the bias is higher with the Canon camera than the Sony camera (R<sup>2</sup> = 0.48 and 0.60, respectively; Bias  
 483 = 43.2 and 33.7, respectively; Figure 13.C); nevertheless, both show similar rRMSE (15% and 16%,  
 484 respectively; Figure 13.C). Figure 13C depicted outliers for both cameras but it is not possible to  
 485 attribute them to one of the cameras or a human error.

486 Overall, the *DeepCount* algorithm showed a solid performance in identifying wheat spikes at early or  
 487 later growth stages. Visual inspection of results also showed that the proposed CNN model was able  
 488 to discriminate ears and background (soil, leaves, etc.) and classified them on a pixel level. The  
 489 proposed model was capable of minimising effects related to brightness, shadow, ear size and shape,  
 490 awn or awnless cultivars and even overlap ears in most scenarios. It should be highlighted that the  
 491 strength of the algorithm also resides in its training data set, where images were collected by a third  
 492 party on completely independent trials, different spatial resolutions, and different varieties than the  
 493 wheat materials in this study. An improvement in the performance would be expected via the  
 494 optimisation of data acquisition process both in the field and through-within images. We believe that  
 495 the optimum configuration is to take images at 2.0-2.5 m above canopies using the focal length between  
 496 35-60 mm which is similar to what human eyes see. Moreover, we noticed that the textural information  
 497 will fade away when spatial resolution is below 0.2-0.3 mm, which will degrade the identification  
 498 performances. ~~It should be highlighted that the strength of the algorithm also resides in its training data~~  
 499 ~~set where images were collected by a third party on completely independent trials, different spatial~~  
 500 ~~resolutions, and different varieties than the wheat materials in this study.~~

### 501 3.3 The effect of nitrogen rate on the performance of the *DeepCount* model

502 We also investigated the effect of nitrogen on the performance of the *DeepCount* method. It was  
 503 expected that the performance of the algorithm declines with the increase of nitrogen use since the  
 504 canopies with a higher level of nitrogen have higher ear density which ears are more overlapped and  
 505 clustered; however, the results showed otherwise. As depicted in Table 2, the overall N3 and N4 data  
 506 had a lower R<sup>2</sup> (0.53 and 0.60, respectively) compared to the overall N1 and N2 data (0.81 and 0.69,

507 respectively). ~~However~~On the other hand, the 2016-FS and 2016-WGIN trials do not follow the same  
 508 pattern. For instance, in the 2016-FS trial, N3 had the highest  $R^2$  value ( $R^2 = 0.89$ ), followed by N2 and  
 509 N1 ( $R^2 = 0.75$  and  $0.59$ , respectively), whereas in the 2016-WGIN, the N4 treatment had the highest  
 510  $R^2$  ( $0.63$ ). Furthermore, on closer inspection, ~~of~~ the N3 and N4 treatments ~~in the 2015-FS, 2016-WGIN,~~  
 511 ~~and combined datasets~~ showed the highest bias values and underestimation of the ear density in the  
 512 2015-FS, 2016-WGIN, and combined datasets. ~~in the automatic method as opposed to the ground~~  
 513 ~~measurements.~~

514 Despite that, the accuracy of the overall experiments for each nitrogen treatment did not change too  
 515 much as the rRMSE value for N1, N2, N3 and N4 were 18, 13, 16 and 15%, respectively. In the end,  
 516 the results did not suggest that the performance of the *DeepCount* model degrades due to the complex  
 517 canopies with a high level of ear density.

## 518 4 Conclusion

519 In this study, the main objective was to present an automatic model that quantifies the number of wheat  
 520 ears in an image or image series. Regardless, of the challenges posed by the acquisition protocol or  
 521 environmental variations in the field, the model was able to deliver the total number of wheat ears  
 522 within an image and/or estimated the number of ears/m<sup>2</sup> if a ground standard was present in the image.  
 523 We demonstrated the feasibility of the proposed technique in which the model was validated on  
 524 numerous images taken from a broad range of spatial resolution images and various data acquisition  
 525 systems. It has been shown that the model can be an essential tool for high throughput analysis and has  
 526 the potential to reduce labour involvement considerably. To minimise the uncertainties between the  
 527 automatic methods and the ground-based measurements, we recommend to 1) have the same sample  
 528 areas 2) have a more reliable ground standard rather than a A4 sheet used in this study 3) take sampling  
 529 from larger area for both image sampling and field measurements 4) increase the spatial resolution of  
 530 visible image to avoid losing the textural information 5) use the focal length of lens between 35– 60  
 531 mm. The code can be found at <https://github.com/pouriaast>

532 In the end, the aim is to increase the adoption of the approach by farmers and breeders by lowering the  
 533 expense of camera equipment. The proposed model can be used as a high-throughput post- processing  
 534 method to quantify the number of spikes for large-scale breeding programs. Furthermore, the automatic  
 535 technique can facilitate farmers to make improved yield estimates, which can be used to plan  
 536 requirements for grain harvest, transport and storage. Subsequently, iImproved estimates could reduce  
 537 post-farm gate costs.

538 The *DeepCount* model benefitted from the CNN architecture and even though the model was trained  
 539 to distinguish two classes, nothing prevents modifying the network to classify and segment more plants  
 540 or species. Given an adequate training model, the proposed semantic segmentation technique offers the  
 541 advantages of versatility and may be applied to other types of applications such as segmenting different  
 542 part of plants organs, vegetation and even detect diseases. In future work, we aim to envisage the use  
 543 of thermal and hyperspectral images which will offer additional information to RGB visible images.

## 544 5 Abbreviations

545 FS Field Scanalyzer

546 CNN Convolutional Neural Network

547 DNN Deep Neural Network

548	NN	Neural Network
549	SLIC	Simple Linear Iterative Clustering
550	WGIN	Wheat Genetic Improvement Network

## 551 **6 Conflict of Interest**

552 The authors declare that the research was conducted in the absence of any commercial or financial  
553 relationships that could be construed as a potential conflict of interest.

## 554 **7 Author Contributions**

555 P.S.T proposed and developed the computer vision methods. P.S.T conducted the image processing  
556 analysis. N.V performed the statistical analysis. N.V planned and conducted the field experiments  
557 under the Scanalyzer. M.J.H contributed to the revision of the manuscript and supervised the project.  
558 All authors gave final approval for publication.

## 559 **8 Funding**

560 Rothamsted Research receives support from the Biotechnology and Biological Sciences Research  
561 Council(BBSRC) of the UK as part of the Designing Future Wheat  
562 ([BB/P016855/1BBS/E/C/00010220](https://doi.org/10.1002/1365-3113.tbr201700010)) and Defra Wheat Genetic Improvement Network (WGIN)  
563 (CH1090) projects.

## 564 **9 Acknowledgments**

565 Authors would like to thank Andrew Riche, David Steele, and March Castle for collecting images from  
566 the WGIN trial. We also thank our colleagues at Limagrain Europe who provided the training datasets.

## 567 **10 References**

- 568  
569 Achanta R, Shaji A, Smith K, Lucchi A, Fua P. Slic superpixels 2010.  
570  
571 Achanta R, Shaji A, Smith K, Lucchi A, Fua P, Süsstrunk S. SLIC Superpixels Compared to State-  
572 of-the-Art Superpixel Methods. IEEE Transactions on Pattern Analysis and Machine Intelligence  
573 2012;34:2274–2282. doi:10.1109/tpami.2012.120 .  
574  
575 Achanta R, Süsstrunk S. Superpixels and Polygons Using Simple Non-iterative Clustering  
576 2017:4895–904. doi:10.1109/cvpr.2017.520 .  
577  
578 den Bergh M, Boix X, Roig G, Gool L. SEEDS: Superpixels Extracted Via Energy-Driven Sampling.  
579 Int J Comput Vision 2015;111:298–314. doi:10.1007/s11263-014-0744-2 .  
580  
581 Busemeyer L, Mentrup D, Möller K, Wunder E, Alheit K, Hahn V, et al. BreedVision — A Multi-  
582 Sensor Platform for Non-Destructive Field-Based Phenotyping in Plant Breeding. Sensors  
583 2013;13:2830–2847. doi:10.3390/s130302830 .  
584  
585 Cointault F, Gouton P. Texture Or Color Analysis In Agronomic Images For Wheat Ear Counting.

- 586 2007 Third International IEEE Conference on Signal-Image Technologies and Internet-Based System  
587 SITIS 2007:696 701. doi:10.1109/sitis.2007.80 .  
588
- 589 Cointault F, Guerin D, Guillemain J, Chopinet B. In-field Triticum aestivum ear counting using  
590 colour-texture image analysis. New Zeal J Crop Hort 2008a;36:117–30.  
591 doi:10.1080/01140670809510227 .  
592
- 593 Cointault F, Journaux L, Miteran J. Improvements of image processing for wheat ear counting.  
594 OrbiUlgBe 2008b.  
595
- 596 Cooley J, Tukey J. An algorithm for the machine calculation of complex Fourier series. AmsOrg  
597 1965.  
598
- 599 Deng J, Dong W, Socher R, Li L, and LK, 2009. Imagenet: A large-scale hierarchical image  
600 database. IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2005  
601 CVPR 2005 n.d. doi:10.1109/cvpr.2009.5206848", "publicationtitle": "2009" .  
602
- 603 Dutta A, Gupta A, Zissermann A. VGG Image Annotator (VIA) n.d.  
604
- 605 Felzenszwalb PF, Huttenlocher DP. Efficient Graph-Based Image Segmentation. Int J Comput Vision  
606 2004;59:167–81. doi:10.1023/b:visi.0000022288.19776.77 .  
607
- 608 Fernandez-Gallego JA, Kefauver SC, Gutiérrez N, Nieto-Taladriz M, Araus J. Wheat ear counting in-  
609 field conditions: high throughput and low-cost approach using RGB images. Plant Methods  
610 2018;14:22. doi:10.1186/s13007-018-0289-4 .  
611
- 612 Harrington P. Machine Learning in Action 2012;5:384.  
613
- 614 Iglovikov V, Shvets A. TernaNet: U-Net with VGG11 Encoder Pre-Trained on ImageNet for  
615 Image Segmentation 2018.  
616
- 617 Jansen M, Dornbusch T, Paulus S, Niehaus B, Sadeghi-Tehran P, Virlet N, et al. Field Scanalyzer –  
618 high precision phenotyping of field crops n.d.  
619
- 620 Kaufman H, Bar-Kana I, Sobel K. Direct Adaptive Control Algorithms: Theory and Applications.  
621 Springer-Verlag 1994.  
622
- 623 Kirchgessner N, Liebisch F, Yu K, Pfeifer J, Friedli M, Hund A, et al. The ETH field phenotyping  
624 platform FIP: a cable-suspended multi-sensor system. Functional Plant Biology 2017;44:154.  
625 doi:10.1071/fp16165 .  
626
- 627 Krizhevsky A, Sutskever I, Hinton G. Imagenet classification with deep convolutional neural  
628 networks. PapersNipsCc 2012:1097 1105.  
629
- 630 Li Z, Chen J. Superpixel Segmentation Using Linear Spectral Clustering. 2015 Ieee Conf Comput  
631 Vis Pattern Recognit Cvpr 2015:1356–63. doi:10.1109/cvpr.2015.7298741 .  
632
- 633 Lomonaco V. Deep learning for computer vision: a comparison between convolutional neural  
634 networks and hierarchical temporal memories on object recognition tasks 2015.



- 635  
636 Madec S, Jin X, Lu H, Solan B, Liu S, Duyme F, et al. Ear density estimation from high resolution  
637 RGB imagery using deep learning technique. *Agr Forest Meteorol* 2019;264:225–34.  
638 doi:10.1016/j.agrformet.2018.10.013 .  
639
- 640 Marsland S. *Machine Learning: An Algorithmic Perspective* 2009.  
641
- 642 Mikolov T, Sutskever I, Chen K, Corrado G. Distributed representations of words and phrases and  
643 their compositionality. *PapersNipsCc* 2013:3111–3119.  
644
- 645 Mohanty SP, Hughes DP, Salathé M. Using Deep Learning for Image-Based Plant Disease  
646 Detection. *Front Plant Sci* 2016;7:1419. doi:10.3389/fpls.2016.01419 .
- 647 Pask AJ, Pietragalla J, Mullan DM, Reynolds MP. *Physiological breeding II: a field guide to wheat*  
648 *phenotyping*. Cimmyt; 2012.  
649
- 650 Pound M, Atkinson J, Wells D. Deep learning for multi-task plant phenotyping.  
651 *OpenaccessThecvfCom* 2017.  
652
- 653 Rahman S, Rahman M, Abdullah-Al-Wadud M, Al-Quaderi G, Shoyaib M. An adaptive gamma  
654 correction for image enhancement. *Eurasip J Image Vide* 2016;2016:35. doi:10.1186/s13640-016-  
655 0138-1 .  
656
- 657 Ren X, Malik J. Learning a classification model for segmentation. *IEEE Computer Society*  
658 *Conference on Computer Vision and Pattern Recognition, 2005 CVPR 2005* 2003.  
659
- 660 Ronneberger O, Fischer P, Brox T. *U-Net: Convolutional Networks for Biomedical Image*  
661 *Segmentation* 2015.  
662
- 663 Schmidhuber J. Deep learning in neural networks: An overview. *Neural Networks* 2015;61:85–117.  
664 doi:10.1016/j.neunet.2014.09.003 .  
665
- 666 Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image recognition.  
667 *arXiv preprint arXiv:1409.1556*. 2014 Sep 4.  
668
- 669 Tu W-C, Liu M-, Jampani V, Surr D, Chien S-Y, Yang M-H, et al. Learning Superpixels with  
670 Segmentation-Aware Affinity Loss. *OpenaccessThecvfCom* 2018:568–76.  
671 doi:10.1109/cvpr.2018.00066 .  
672
- 673 Ubbens J, Cieslak M, Prusinkiewicz P, Stavness I. The use of plant models in deep learning: an  
674 application to leaf counting in rosette plants. *Plant Methods* 2018;14:6. doi:10.1186/s13007-018-  
675 0273-z .  
676
- 677 Virlet N, Sabermanesh K, Sadeghi-Tehran P, Hawkesford MJ. Field Scanalyzer: An automated  
678 robotic field phenotyping platform for detailed crop monitoring. *Funct Plant Biol* 2016;44:143–53.  
679 doi:10.1071/fp16163 .  
680
- 681 Wang M, Liu X, Gao Y, Ma X, Soomro NQ. Superpixel segmentation: A benchmark. *Signal Process*  
682 *Image Commun* 2017;56:28–39. doi:10.1016/j.image.2017.04.007 .

683  
 684 Xiong X, Duan L, Liu L, Tu H, Yang P, Wu D, et al. Panicle-SEG: a robust image segmentation  
 685 method for rice panicles in the field based on deep learning and superpixel optimization. Plant  
 686 Methods 2017;13:104. doi:10.1186/s13007-017-0254-7 .

687  
 688 Zhou C, Liang D, Yang X, Yang H, Yue J, Yang G. Wheat Ears Counting in Field Conditions Based  
 689 on Multi-Feature Optimization and TWSVM. Front Plant Sci 2018;9:1024.  
 690 doi:10.3389/fpls.2018.01024 .

691

692

693

694

## Figures

695 **Figure 1 Schematic representation of the *DeepCount* method**

696

697 **Figure 2 Over-head view digital images of wheat cultivars with different canopy complexity**  
 698 **taken in the field using the handheld DSLR camera (A and B) and the Field Scanalyzer platform**  
 699 **(C). An A4 sheet is placed over the canopy for each image as a ground standard system to**  
 700 **transform the total number of wheat ears in the image into number of ears/m<sup>2</sup>**

701

702 **Figure 3 Training patches. Examples of expert annotation of spikes for different wheat cultivars**  
 703 **without awns (A), with awns (B), and backgrounds (e.g. soil, leaves)**

704

705 **Figure 4 Augmented samples of the same spike with various transformations such as random**  
 706 **zoom, rotation, flipping, brightness and gamma correction. For example, 1) the original image;**  
 707 **5 & 10) adjusted HSV colour image; 6, 8 & 10) gamma colour correction. Cropping, flipping,**  
 708 **zooming and scaling was applied to all image randomly with the probability of 0.5.**

709

710 **Figure 5 Examples of superpixel segmentation using the SLIC technique.**

711

712 **Figure 6 A.1 & B.1 show the SLIC superpixel outputs. A.2 & B.2 are the results of pixel-wise**  
 713 **semantic segmentations. The red circle illustrates the imperfection in the SLIC method.**

714

715 **Figure 7 A.1 & B.1 show the SLIC superpixel outputs. A.2 & B.2 are the output of the**  
716 ***Deepcount* model. The red circle illustrates the imperfection in the SLIC method**

717

718 **Figure 8 Encoder-decoder neural network architecture also known as U-Net where VGG-16**  
719 **neural network without fully connected layers as its encoder. The number of channels increase**  
720 **stage by stage on the left part while decrease stage by stage on the right decoding part. The**  
721 **arrows show transfer of information from each encoding layer and concatenating it to a**  
722 **corresponding decoding part**

723

724 **Figure 9 A plot of loss and accuracy over the course of 15 epochs with a 1e-4 learning rate. Using**  
725 **of pre-trained VGG model on ImageNet dataset helped the model to converge quicker**

726

727 **Figure 10 The hand-crafted ear-counting method. A) original image B) greyscale image C)**  
728 **result after applying edge detection technique D) dilate the image E) fill the holes F) filtering by**  
729 **removing small objects (noises) G) erode and smooth the image H) counting the contours/ears**

730

731 **Figure 11 Examples of result images A) WGIN experiment with an A4 sheet used as a ground**  
732 **standard B) Field Scanalyzer experiment in 2015**

733

734 **Figure 12 Comparison of the number of ears visually annotated on the images (Annotation – A,**  
735 **B) and the number of ears/m<sup>2</sup> (C, D) with the number of ears estimated by the Edge (A, C) and**  
736 ***DeepCount* (B,D) methods for the 2 dataset collected with the Field Scanalyzer in 2015 (blue dots)**  
737 **and 2016 (red triangles)**

738

739 **Figure 13 Comparison between the number of ears/m<sup>2</sup> counting form the field and the number**  
740 **of ears estimated by the ~~neural-network~~DeepCount model (NN) method for the datasets collected**  
741 **with the Field Scanalyzer in 2015 (A – open circle) and in 2016 (B - open triangles), for the WGIN**  
742 **trial in 2016 (C – cross) separated by camera (D), for all datasets together (E) and for all dataset**  
743 **together separated by nitrogen level (f - N1: blue, N2: green, N3: red and N4: purple)**

744

## Tables

745 **Table 1 Characteristics of the three experiments considered in this study**

Dataset	Plot	Nitrogen (kg/ha)	Image	Camera	Image size	Focal length	Resolution (mm)	Date
2015 - FS	72	0, 100, 200, 350	72	Prosilica GT 3300 Allied Vision	3296×2474	50 mm	0.22-0.29	13/07/2015
2016 - FS	54	0, 100, 200	54	Prosilica GT 3300 Allied Vision	3296×2474	50 mm	0.26	29/06/2016
			78	Canon G12	3648×2736	6 mm	0.21-0.31	13/06/2016
2016-WGIN	499360	0, 100, 200, 350	121	SONY - NEX-7	6000×3376	18 mm	0.14-0.25	13/06/2016

746

747 **Table 2 Comparison between the number of ears/m<sup>2</sup> counting form the field and the number of**  
748 **ears estimated by the *DeepCount* model for the 3 datasets collected, separately and combined**  
749 **for each of the nitrogen levels. Performance of the *DeepCount* model in regard to the nitrogen rate**  
750 **across the different experiments- a and b are the slope are the offset of the regression line, respectively.**

		N1	N2	N3	N4
2015 - FS	a	1.16	0.96	0.64	0.68
	b	-18.40	55.22	263.06	282.56
	R <sup>2</sup>	0.58	0.46	0.15	0.22
	RMSE	61.50	60.30	92.20	122.90
	rRMSE	13%	10%	14%	17%
	Bias	42.30	35.20	58.90	100.10
2016 - FS	a	0.75	1.10	0.93	
	b	45.23	-16.13	39.29	
	R <sup>2</sup>	0.59	0.75	0.89	
	RMSE	41.00	43.80	32.40	
	rRMSE	22%	10%	7%	
	Bias	-19.60	20.20	9.70	
2016 - WGIN	a		0.95	0.71	0.88
	b		33.87	189.27	89.62
	R <sup>2</sup>		0.42	0.41	0.63
	RMSE		67.10	98.30	72.00
	rRMSE		15%	17%	14%
	Bias		15.60	57.90	30.80
All dataset	a	1.28	1.06	0.83	0.96

b	-76.43	-4.10	131.26	66.39
R <sup>2</sup>	0.81	0.69	0.53	0.60
RMSE	52.20	61.40	90.00	84.40
rRMSE	18%	13%	16%	15%
Bias	11.30	20.70	50.30	44.30

---

751

In review

Figure 1.JPEG

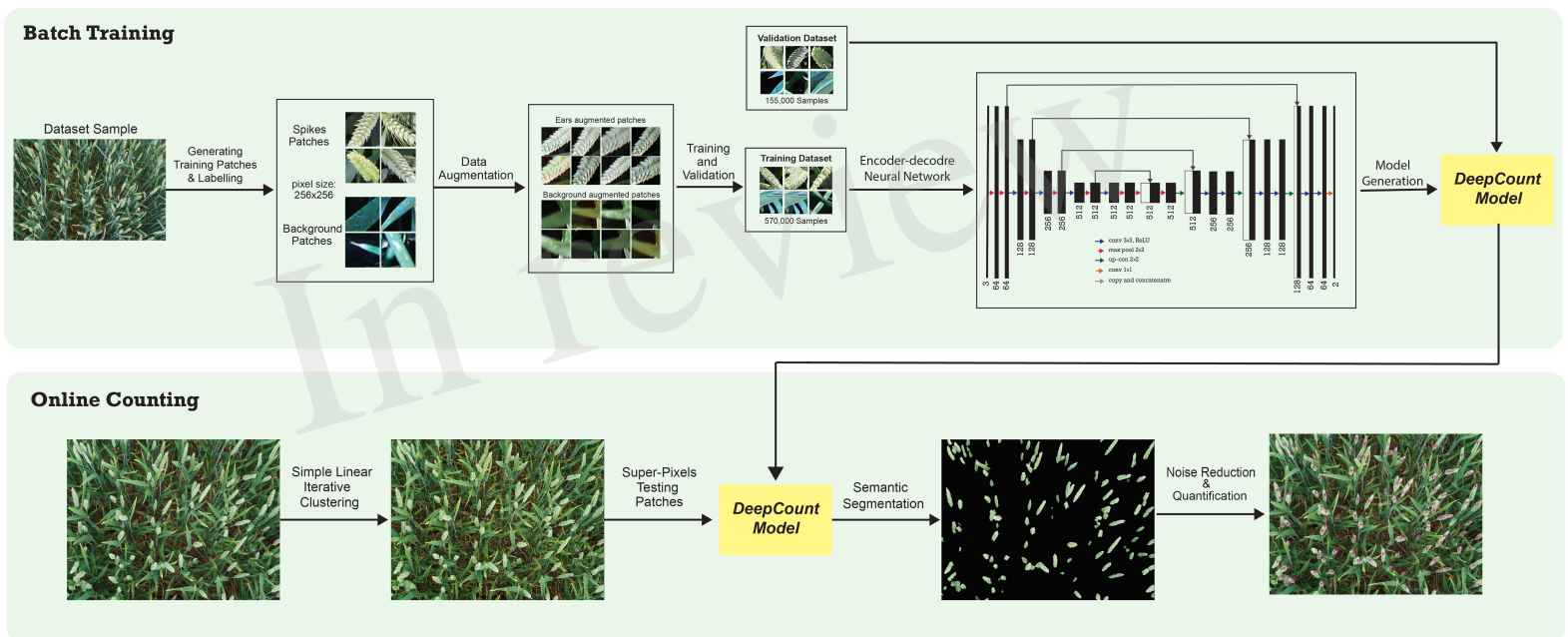




Figure 2.JPEG





Figure 3.JPEG

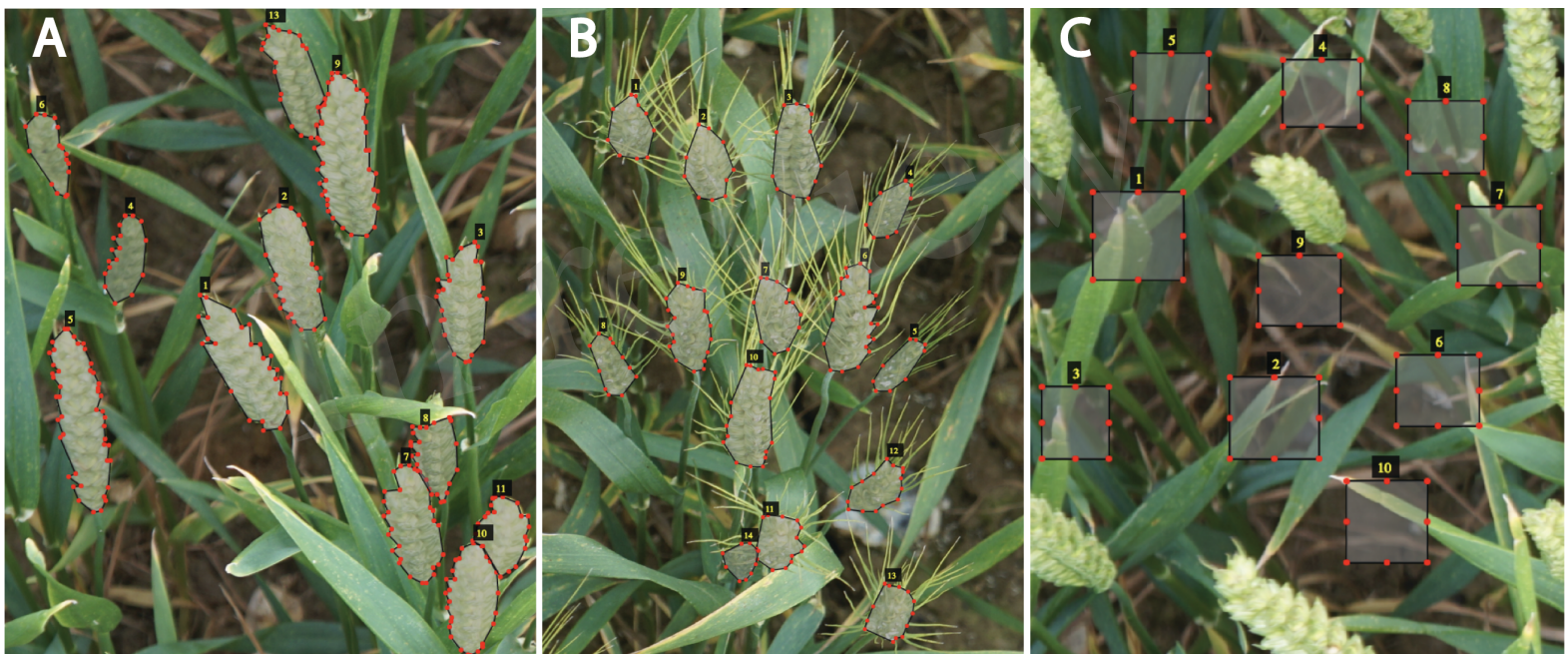




Figure 4.JPEG

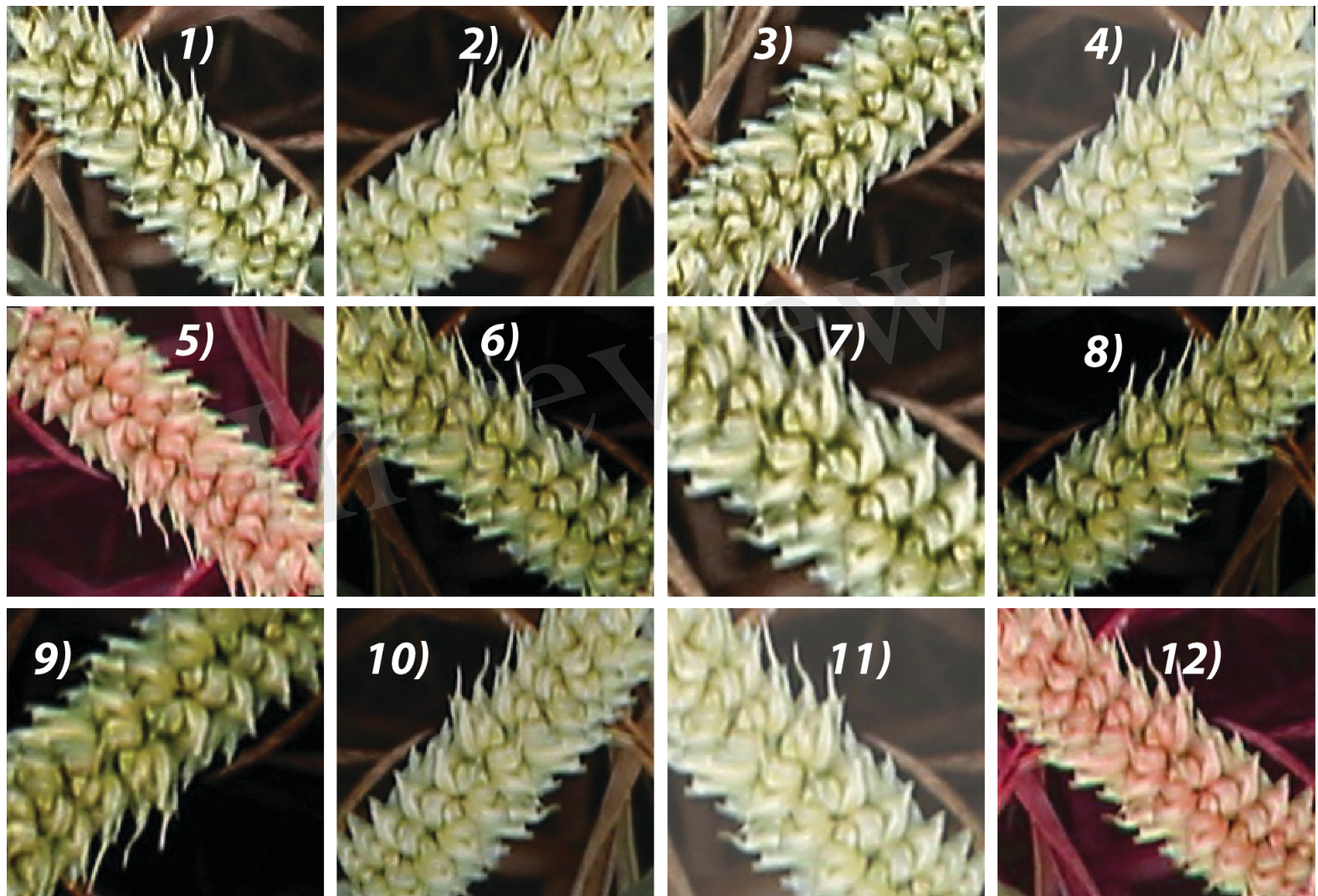


Figure 5.JPEG

In review

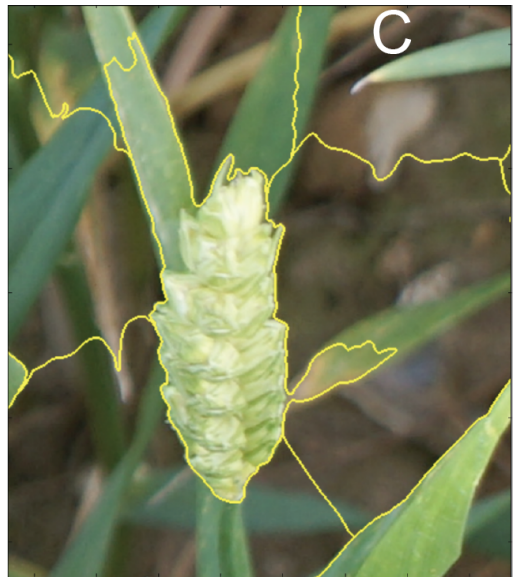
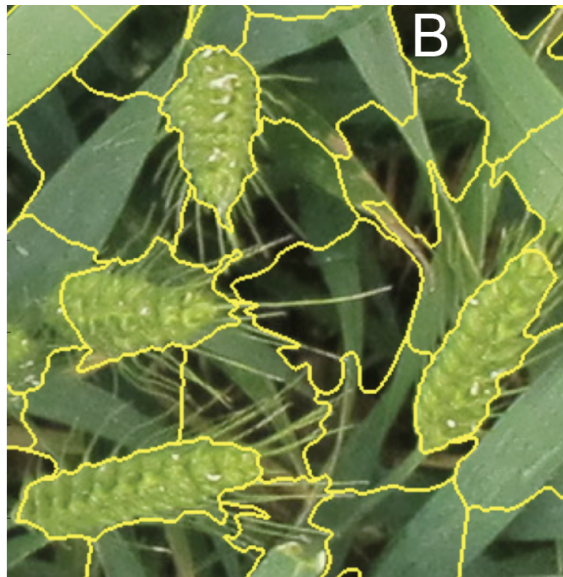




Figure 6.JPEG

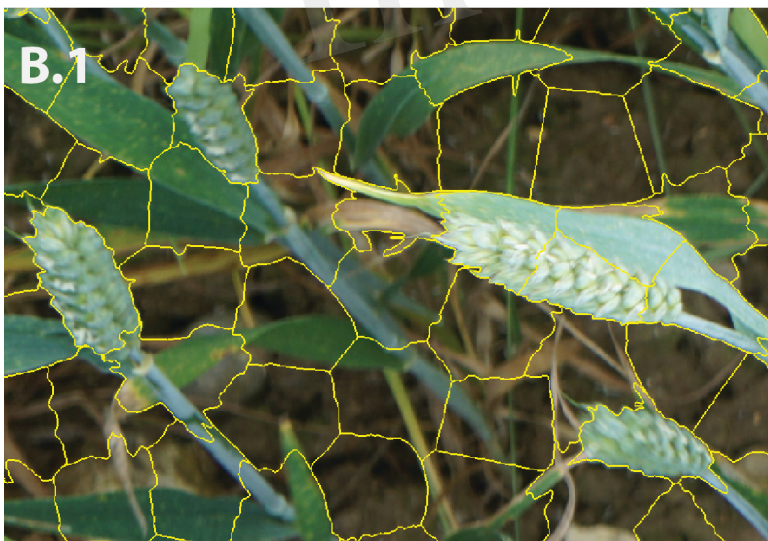


Figure 7.JPEG



Figure 8.JPEG

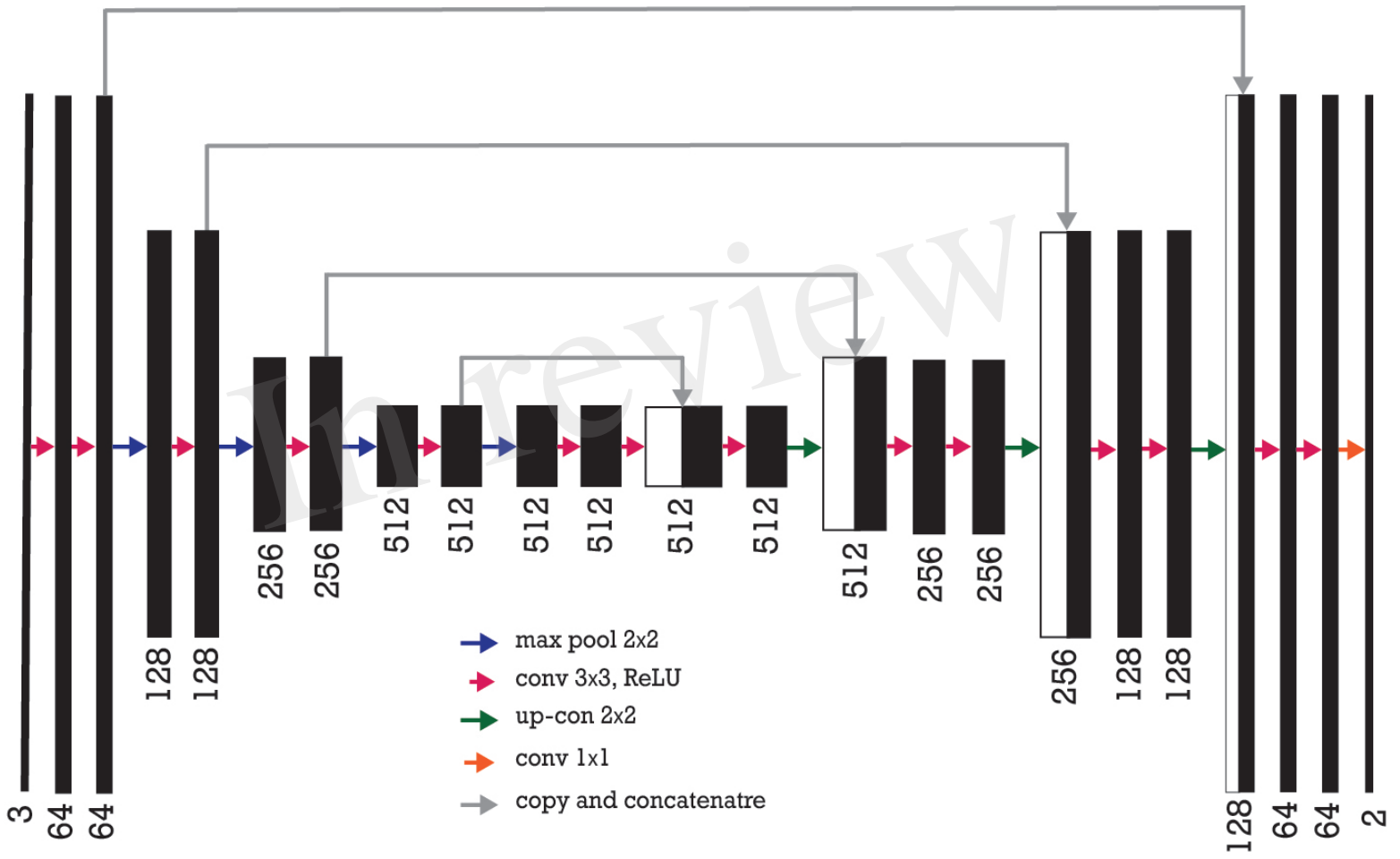


Figure 9.JPEG

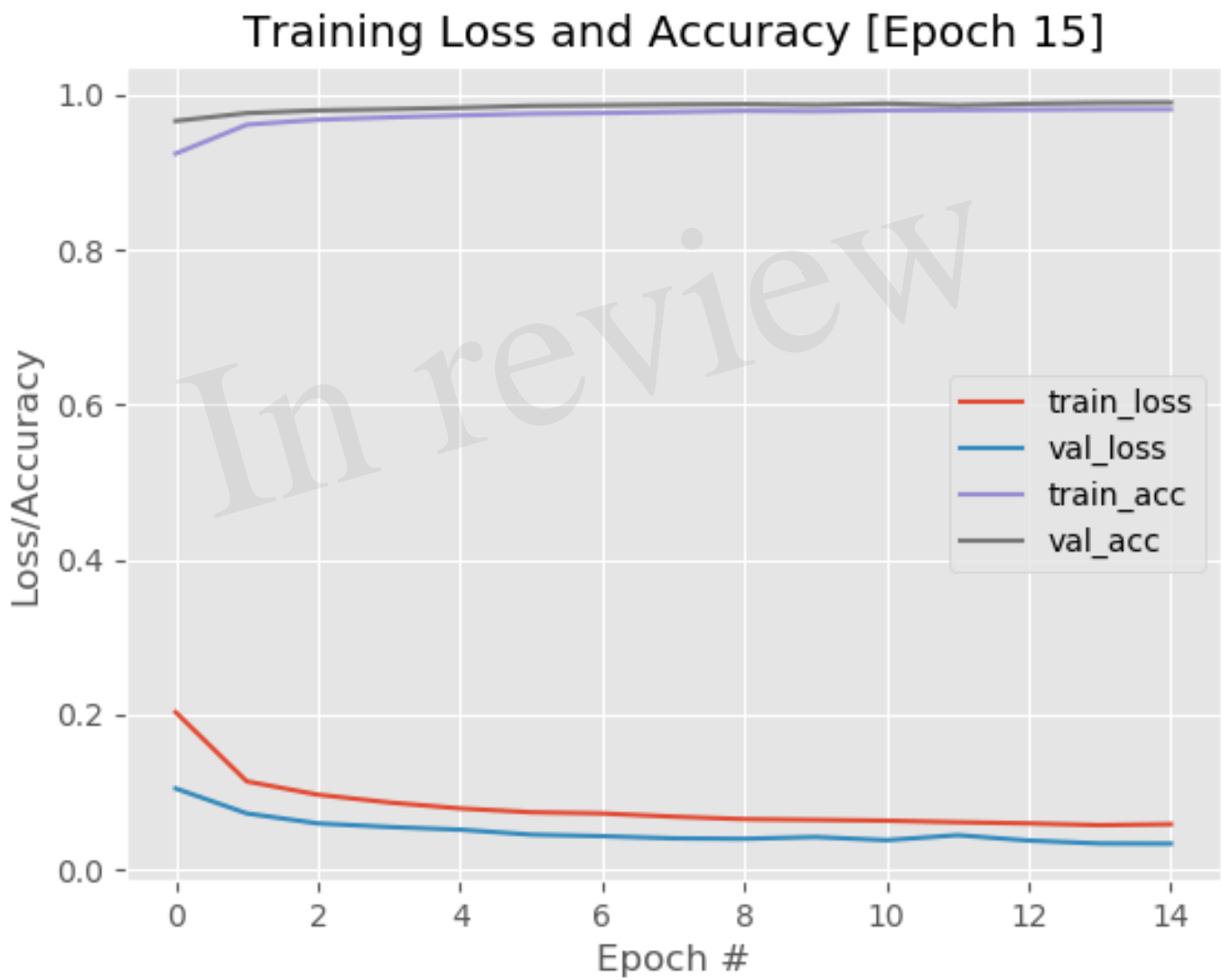


Figure 10.JPEG

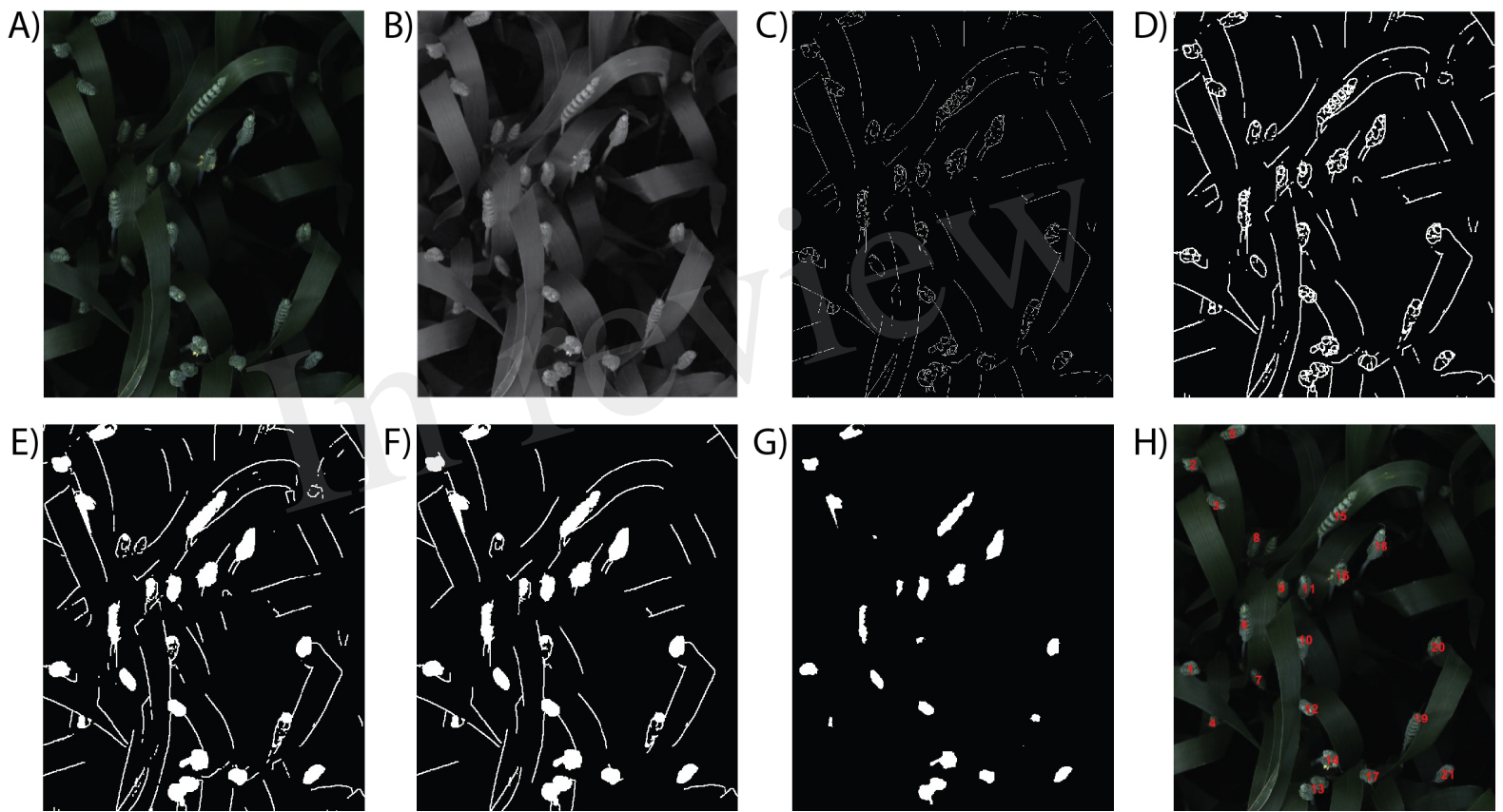




Figure 11.JPEG

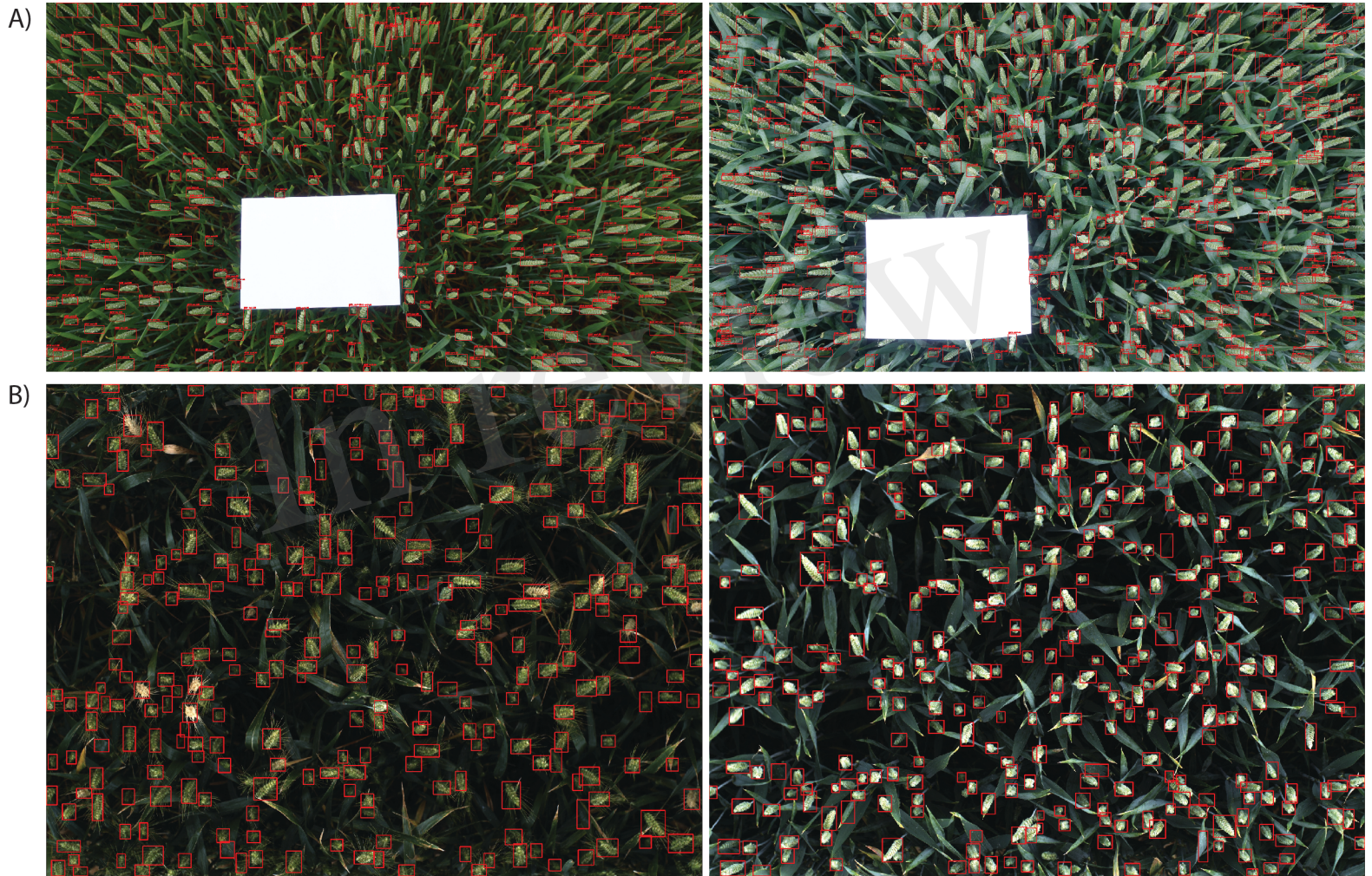




Figure 12.JPEG

