



HAL
open science

Novel data augmentation strategies to boost supervised segmentation of plant disease

Clément Douarre, Carlos F Crispim-Junior, Anthony Gelibert, Laure Tougne,
David Rousseau

► To cite this version:

Clément Douarre, Carlos F Crispim-Junior, Anthony Gelibert, Laure Tougne, David Rousseau. Novel data augmentation strategies to boost supervised segmentation of plant disease. *Computers and Electronics in Agriculture*, 2019, 165, pp.104967. 10.1016/j.compag.2019.104967 . hal-02330900

HAL Id: hal-02330900

<https://hal.science/hal-02330900v1>

Submitted on 20 Jul 2022

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

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



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

Novel data augmentation strategies to boost supervised segmentation of plant disease images

Clément Douarre^{a,b,c,*}, Carlos F. Crispim-Junior^a, Anthony Gelibert^c, Laure Tougne^a, David Rousseau^b

^aUniv Lyon, Université Lyon 2, LIRIS, UMR CNRS 5205, F-69676, Lyon, France

^bLaris, UMR INRA IRHS, Université d'Angers, 62 Avenue Notre Dame du Lac, 49000 Angers, France

^cCarbon Bee, Rue du Commerce, ZA Les Plaines, 26320 St Marcel-Lès-Valence, France.

Abstract

Annotation of images in supervised learning is notably costly and time-consuming. In order to reduce this cost, our objective was to generate images from a small dataset of annotated images, and then use those synthesized images to help the network's training process. In this article, we tackled for illustration with agricultural material the difficult segmentation task of apple scab on images of apple plant canopy by using convolutional neural networks. We devised two novel methods of generating data for this use case: **one based on a plant canopy simulation and the other on Generative Adversarial Networks (GANs)**. As a result, we found that simulated data could provide an important increase in segmentation performance, up to a 17% increase of F1 score (a measure taking into account precision and recall), compared to segmenting with weights initialized on ImageNet. In this way, we managed to obtain, with small datasets, **higher segmentation scores** than the ones obtained with bigger datasets if using no such augmentations. Moreover, we left our annotated dataset of scab available for the plant science imaging community. The proposed method is of large applicability for plant diseases observed at a canopy scale.

1. Introduction

In computer vision applied to agriculture, machine learning techniques are currently progressing very rapidly (Minervini et al., 2015; Pound et al., 2017; Kamilaris and Prenafeta-Boldú, 2018). Two categories of methods are usually distinguished in machine learning depending on whether it is supervised or unsupervised. Supervised techniques require annotation, a task that is possibly time-consuming when dealing with images with complex structures, such as texture, or with large dataset of images, as in phenomics for plant sciences. This task is even more time-consuming now with the success of deep learning methods. These methods show great success (Goodfellow et al., 2016) on a variety of image processing tasks of seemingly unlimited complexity thanks to sufficiently large neural networks. However, as these networks are composed of many computational units, they typically require a great quantity of annotated data compared to other learning methods. Yet, in the case of plant disease detection, annotated datasets of infected plants are often very small. Often, the severity of a given infection is judged by the quantity of visible lesions: Hence, image segmentation *i.e.* pixelwise classification is needed, which requires pixelwise annotation, which is tedious and difficult.

Our goal was to obtain the best segmentation score possible using a neural network based segmentation and data simulation techniques with very small size of the datasets. We investigated the possibility to circumvent manual annotation by the generation of synthetic annotated data. Such generators fall in three broad categories: (i) Standard data augmentation, which has just started to be tested in plant sciences, (ii) computer graphics techniques of image rendering and (iii) generative adversarial networks (GANs) (Goodfellow et al., 2014) which is a new approach, not based on computer graphic model but on machine learning. These three kinds of generators have only recently been tested separately in plant sciences (Pawara et al., 2017; Ubbens et al., 2018; Giuffrida et al., 2017) mainly with images of single healthy plants. In this work, we implemented and compared an example of each of these categories on the same plant disease use case at the observation scale of the canopy.

We focus on apple scab which is an infection afflicting apple trees and in general trees of the *Malus* genus. This infection is caused by fungus *Venturia inaequalis*, which leads to brown lesions on the leaves and the fruits. It requires more than ten fungicide treatments per year to be controlled and can be considered as the most serious disease for apple plantations (Bowen et al., 2011). Quantifying the development of scab is of great importance for studying interactions of apple scab and apple trees, as well as for analyzing the evolution of pathogenicity in *Venturia inaequalis* populations and for breeding scab-resistant apple cultivars. In order to improve its treatment,

*Corresponding author

Email addresses: clement.douarre@liris.cnrs.fr (Clément Douarre), carlos.crispimjunior@univ-lyon2.fr (Carlos F. Crispim-Junior), anthony.gelibert@carbonbee.fr (Anthony Gelibert), laure.tougne@liris.cnrs.fr (Laure Tougne), david.rousseau@univ-angers.fr (David Rousseau)

early detection, *i.e.* before visible symptoms would be valuable. The life cycle of the fungus starts by a phase under the leaf’s cuticle (invisible to the naked eye), before rupturing to the leaf’s surface, at which point infection is visible but already very advanced (Oerke et al., 2011). However, the early stage is detectable by observation of other wavelengths than the visible spectrum. Indeed, in that stage, the fungus blocks the host leaf’s water evacuation, which in turn modifies the temperature of infected zones. Infrared imagery, which provides information on temperature, is therefore suitable for early detection (Belin et al., 2013; Delalieux et al., 2009). The infrared domain has shown optimal contrast compared to RGB (Oerke et al., 2011; Chéné et al., 2012; Belin et al., 2013; Benoit et al., 2016). The detection of apple scab was implemented on individual plants at a very early stage of development where plants had a very limited amount of leaves (Oerke et al., 2011; Chéné et al., 2012). In a real agronomic and biological research environment, populations of trees are grown and the observation scale of plants is rather at canopy level.

This work is positioned along the following lines toward the closest related literature on data augmentation. Most simulation of plant images has been conducted by the implementation of biological models to generate 3D scenes. These simulations are then coupled with an artificial acquisition system in order to generate 2D images. Full plant simulation framework have been developed, enabling researchers to model a wide variety of crops and plants, along with many physical measurements on these models (Pradal et al., 2008). Some of these simulations were used to boost machine learning algorithms. For example simulator of root images were developed, which also generated annotation masks for a segmentation task (Benoit et al., 2014). Closer to our work, some simulators were designed to serve as dataset augmentation. In recent years, a weed simulator was designed to improve weed detection (Di Cicco et al., 2017), and high-quality 3D scenes of pepper plants were generated for a model of species identification (Barth et al., 2018). Some researchers focused, like this work, at the leaf level: 3D *Arabidopsis* leaves were generated in order to improve leaf segmentation (Ward et al., 2018). While these elaborate models generate highly realistic images, they often require many complex steps going through sometimes proprietary software, such as Blender. The model-based simulator presented here is deliberately simpler than the ones presented. It was simply based on a Python script and an open-source dataset, and worked in two dimensions. Also, it simulated scab disease, which is a finer level of detail than simulating whole species. Generative Adversarial Networks have now much improved from their debuts (Goodfellow et al., 2014). While the original GAN strived to generate images from a given true image distribution, the key concept of generating images for a given conditional distribution was then introduced (Mirza and Osindero, 2014). The condition the authors used was a class label: the conditional GAN, or cGAN, learned to generate images of specific digits from the MNIST dataset (instead of learning to generate all digits indifferently). This was shown to actually improve the visual quality of GAN results on that dataset compared to using no labels, and it also paved the way for a new range of applications for this kind of network. For example, one example of using such cGANs for data synthesis is generating different kinds of liver legion images in order to improve their classification (Frid-Adar et al., 2018). The cGAN’s principle was generalized by providing images and strings as inputs instead of class labels (Isola et al., 2017). Images as inputs contain much more information than integer labels. Thus, the objective of such cGANs shifted from actual generation to domain transform, *e.g.* transforming maps to satellite images. Segmentation was also studied by (Isola et al., 2017), considering the image and its annotation as the two domains of interest. While results were quite visually impressive, such segmentation was not competitive with benchmarks of dedicated segmentation networks. The process of generating annotation from images (or vice-versa) has since been applied as a mean of data synthesis to augment datasets before performing segmentation with those dedicated networks. In plant sciences, a recent example concerns *Arabidopsis* leaf counting, where a cGAN was trained to generate plant images using the number of leaves as the label (Giuffrida et al., 2017). Similar work has been conducted where the cGAN’s labels were masks of manually segmented *Arabidopsis* leaves (Zhu et al., 2018). New segmented masks were then generated and used as labels to the cGAN. The approach for data augmentation by GAN generation we explored differed from these attempts, as we did not wish to have to generate plausible masks to train a cGAN on as in (Zhu et al., 2018). Indeed, this is time-consuming, introduces new parameters to the model and requires expertise that is not always available. The architecture we used was capable of generating images and their annotation simultaneously and had therefore a more generative nature than that of (Isola et al., 2017). Our model was similar to works where GANs were trained on image-annotation pairs (Neff et al., 2017). The experiment for which data was generated in their case was however quite different, as it was mainly showed that for the concerned dataset, replacing part of the training images by simulated images yielded results that were similar to those obtained using only the real images. By contrast here, we showed situations where adding simulated images to the train set could actually improve results.

To sum up, our paper focuses on novel data augmentation strategies applied to boost the supervised segmentation of apple scab observed at a canopy level. No prior work has been done on apple scab segmentation at the canopy level. Thus, we compare segmentation results when training on a "full" training dataset, *i.e.* a training dataset

106 where adding new annotated images does not significantly improve results to *a reduced part of this training dataset*
 107 *augmented by our strategies*. This course of action is illustrated in Fig. 1. The rest of the article is organized as
 108 follows. Our methodology is presented in section 2. Results of the experiments are presented in section. 3 and
 109 discussed in section. 4.

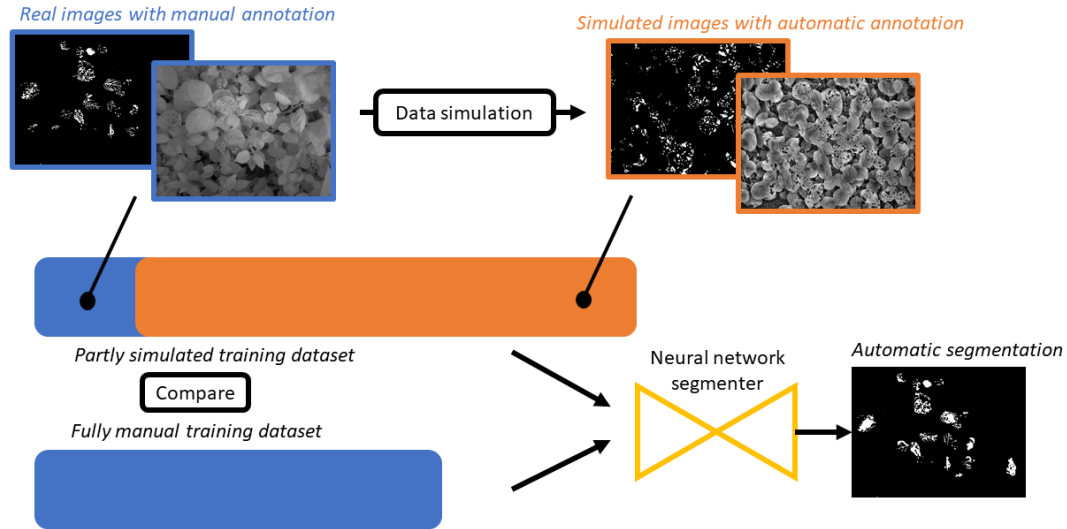


Figure 1: Global illustration of the work. We simulated annotated data from a small set of real annotated images through various simulators. We tested if simulated data increased segmentation score, and also compared this hybrid dataset to a "full" dataset with more real images - whose creation required more manual annotation than for the hybrid dataset.

110 2. Material and methods

111

112 2.1. Data acquisition

113 Images of apple plants were acquired in April 2018 at INRA Angers greenhouses. The camera used was developed
 114 by company Carbon Bee, capable of acquiring images both in visible and infrared light. The camera was held by
 115 hand at 2m above the ground, facing downwards. The apple plants in the camera's field of view were set on trays
 116 and peaked at about 1.5m from the ground. Acquisitions were done around noon and the camera gain was set for an
 117 optimal dynamic range without saturation. A dataset of 11 infrared (IR) images of 2592×1944 pixels, representing
 118 apple (*Malus × domestica Borkh*) plants inoculated with scab (Fig. 2) was acquired. While the number of images
 119 may seem small for a deep learning-based approach, it should be noted that each image represented 50 apple plants,
 120 which led to several hundred visible leaves, among which 15 to 40 were afflicted with scab. We annotated this
 121 dataset using GIMP by assigning either a "scab" or a "not scab" label to each pixel as visible in Fig. 2. This
 122 dataset presented several challenges from a computer vision point of view:

- 123 • The segmentation task at hand concerned small local structures to be localized in an image several orders of
 124 magnitudes larger than them. Hence, structural information, such as texture, was somewhat limited as scab
 125 lesions were captured from an important distance compared to their size.
- 126 • Localization of scabbed leaves was fairly easy for human eyes, but a pixelwise localization was difficult because
 127 of a "gradient" aspect of scab. The amount of fungi on apple scab spots was indeed more concentrated at the
 128 center of these spots and this created a fading of the contrast. Thus, pixelwise annotation was very challenging
 129 even for an agronomic expert, and we expected it to be a challenge for segmentation networks as well.
- 130 • Because of the difference in reflectance that scab lesions cause on leaves, scab could seem easily detectable at
 131 first glance. However, some factors complicated this distinction: other structures such as leaf veins also caused
 132 contrasted areas on the leaves. Moreover, the multi-layer structure of the scene led to leaves that showed very
 133 different gray level intensities depending on their distance from the soil and occlusions by other leaves.

- 134 • Images were unevenly illuminated due to spatial non-homogeneity in natural lighting. Moreover, they suffered
135 from a slight vignetting effect, meaning more light reached the sensor’s center than its borders, leading to images
136 with dark borders.

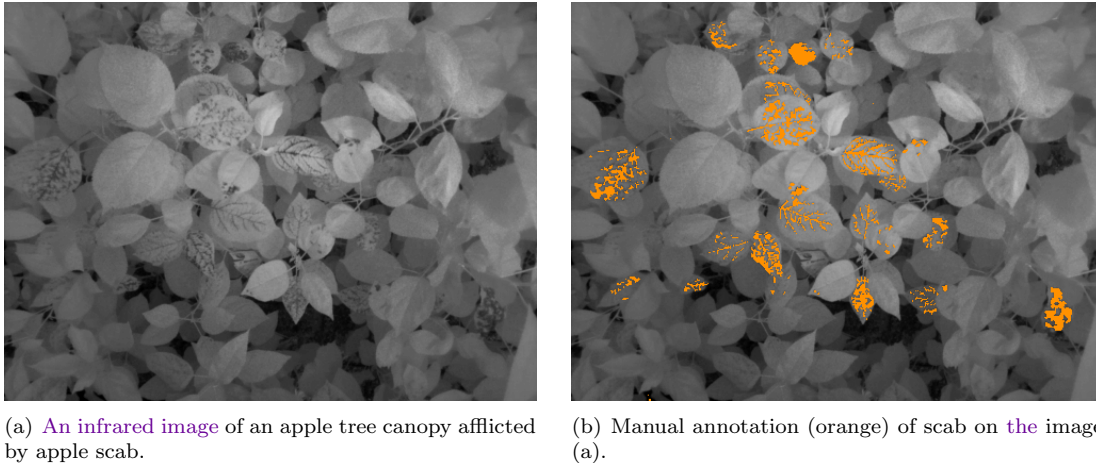


Figure 2: Illustration of plant disease and observation scale considered in this article for the comparative test of data augmentation strategies.

137 2.2. Network training

138 This section presents the segmentation network and training specifics we used for all of our segmentation
139 experiments.

140 2.2.1. Architecture

141 We used the SegNet segmentation network (Badrinarayanan et al., 2017). While not being state-of-the-art for
142 segmentation tasks anymore (Cordts et al., 2018), the goal of this paper was to prove the usefulness of simulated
143 images; we only compared results on different train sets and did not strive for absolute best performance. We
144 therefore chose SegNet because robust implementations were available and its inner workings were extensively
145 studied and understood (Noh et al., 2015) compared to more recent networks. For all experiments, we used
146 the following parameters: learning rate = 10^{-4} , batch size = 5, SGD optimizer with momentum = 0.99 and
147 weight decay = 5×10^{-4} . These parameters were set by starting with standard values for image segmentation with
148 SegNet, then adjusted through a grid search to prevent overfitting and to yield the best results on the validation
149 set. We implemented class balancing because of the high class imbalance in the dataset (Badrinarayanan et al.,
150 2017).

151 2.2.2. Dataset preparation

152 The original images were tiled into sub-images of 64×64 pixels, with no overlap. This was done for computational
153 efficiency and to enable easier generation of images by GAN, which are known to be stable for images of these sizes.
154 During the tiling process, only tiles where at least one pixel had been annotated as "scab" were kept. This strongly
155 helped to reduce class imbalance: in the original dataset, only 2% of pixels were annotated as "scab". After tiling,
156 the scab proportion was 10%. Tiling was not as harmful for segmentation for this dataset as it could have been
157 for other types of images. Images from the dataset represented a global scene of apple plants, while symptoms of
158 interest were fairly local: scab symptoms were independent from plant to plant, and from leaf to leaf. In other
159 words, the general structure of the scene brought almost no information to the local structure of the features of
160 interest. This is not the case in other well-known segmentation datasets such as CamVid (Brostow et al., 2009),
161 where general spatial structure is important. It would therefore have brought no additional information to the
162 network to train on images where several plants were seen together than to train on these plants alone. The dataset
163 was split into train, validation and test sets with respectively 535, 178 and 178 tiles. Out of the 535 tiles of the
164 training set, 107 (20%) were kept as the "reduced dataset" and used in our experiments. The rest of the dataset
165 was kept for comparison of performance with real images. We believed that this reduced quantity of data would be
166 small enough for data augmentation strategies to be useful, and large enough to ensure stable stable segmentation
167 results (in particular regarding overfitting). Moreover, annotating datasets for segmentation is more manageable
168 when these datasets are small. We believe that datasets the size of the "reduced dataset" are close in size to many
169 of the ones used for machine learning in plant sciences.

170 *2.2.3. Training and evaluation*

Training was done on the train set and **stopped** when the score on the validation set had not **been** improved for 50 epochs. The inference was then carried out on the test set with weights that had yielded the lowest validation loss. The well-known metric "accuracy" (true positives + true negatives over whole population) would have been inadequate to assess the quality of a detection on such an unbalanced dataset: for example, a detection of no scab in an image where 10% of an image were scab led to an accuracy of 90%. Therefore, the metric used for performance evaluation was the F1 score (referred to as simply "score" when discussing the results in section 3), defined as

$$F1 = 2 \frac{precision \times recall}{(precision + recall)}.$$

171 This metric is the harmonic mean of precision and recall, defining scab as the class of interest. Precision is the
172 ratio between how many scab pixels the algorithm has marked as such (true positives) and how many pixels the
173 algorithm has marked as scab (true positives + false positives). Recall is the ratio between true positives and how
174 many scab pixels are in the image (true positives + false negatives). A segmentation with a high F1 score therefore
175 has high precision and high recall, meaning few false positives and false negatives. This metric was a relevant choice
176 in such 2-class segmentation cases with heavy class imbalance (Brownlee, 2015).

177 *2.3. Data augmentation strategies*

178 The studied dataset was quite small and segmentation results were quite poor. In this section, we present the
179 different augmentation strategies we implemented for the dataset. Two of them are novel and are described in
180 **details** in Sections. 2.3.2 and 2.3.3. Our implementation of standard data augmentation, which is a type of data
181 simulation very commonly used in neural network based segmentation, is also presented in section. 2.3.1.

182 *2.3.1. Standard data augmentation*

183 When we implemented standard data augmentation (SDA), we considered the following transformations (**Pawara**
184 **et al., 2017**): flip (horizontal or vertical), rotation (90, 180 or 270 degrees), perspective transform by homography,
185 scaling ($\times 0.8$ - $\times 1.2$) and Gaussian blur ($\sigma = 1 - 2$). These transformations were justified by the fact that they
186 could actually occur across images. Rotations and flips were suitable here, as images represented mainly texture
187 and were therefore globally isotropic. Scaling and homography were justified by the fact that images had been
188 acquired with an infrared camera manually held over the apple tree at a distance: the camera could have fluctuated
189 during the acquisition trial. This fluctuation also justified Gaussian blur: the focal of the camera was fixed so that
190 fluctuations of the distance of the camera to plants resulted in a blurring effect. Augmenting a given image meant
191 applying each of these transformations with a certain probability. These probabilities ranged from 0.2 to 0.5. The
192 exact probability depended on the transformation and was meant to reflect the empirical frequency we thought that
193 specific transformation would be useful for the model to generalize better. We devised two ways of augmenting the
194 data:

- 195 • **offline SDA**: We augmented each image in the train set a number n_{aug} of times, multiplying the train set size by
196 n_{aug} .
- 197 • **online SDA**: The train set was left untouched before training starts, but during the training process, transfor-
198 mations were applied to each (original) image at each epoch. It was quite possible that for a relatively small
199 number of epochs, a given image is transformed differently at each epoch. Therefore, the number of different
200 images shown to the network during training with such SDA was actually close to $n_{train} \times n_{epochs}$.

201 *2.3.2. Model-based simulator*

202 We now present **the simulator that** we devised specifically for generating canopy images. It assumed (i) that the
203 scenes of Fig. 2 could be reproduced by replication of leaves put on top of one another similarly to the "dead leaves
204 model" (Matheron et al., 1975) and (ii) that the apple scab on apple leaves could be characterized as a texture
205 defined by its first and second order statistics. The general pipeline was the following:

- 206 1. Randomly draw an isolated apple leaf image from an existing leaf dataset;
- 207 2. Generate simulated scab texture and apply it on the leaf;
- 208 3. Place the "scabbed" leaf in the scene in a structured way.

209 The procedure is illustrated in Fig. 3, and the different steps are explained more thoroughly below.

210

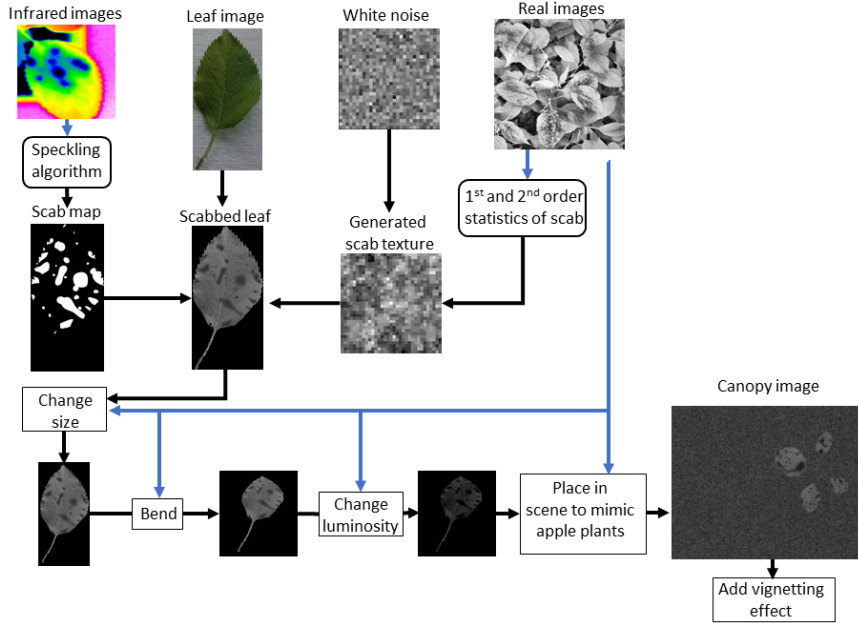


Figure 3: Model-based simulator algorithm. Information from real images and real leaves from a dataset (blue arrows) were used to create a realistic-looking scene. An example output can be seen in Fig. 4.

1. Leaf images: We used the LeafSnap dataset (Kumar et al., 2012) of apple tree leaves set flat and isolated on a uniform background. Images from the train set were used to capture statistics necessary for scab texture synthesis.

Regarding the texture generation step, we chose to implement a "procedural" process, meaning texture was generated from a filtered Gaussian white noise (McCombs, 2005). The process consisted in fitting the first (mean, standard deviation) and second (autocorrelation) order statistics to the ones recorded on the train set. This simple procedure assumed the apple scab contrast is stationary and captured by a Gaussian random process.

2. Scab pattern: To simulate scab placement on leaves, we studied typical shapes of scab infection in infrared light (Oerke et al., 2011). To reproduce this spatial distribution, we used a speckling algorithm that gave results visually similar to real images. This algorithm is based on taking the Fourier transform of a 2D array of 20×20 random complex numbers and yields a pattern of spots whose sizes are controllable. We then added a Gaussian intensity to these "spots" to produce images closer to real ones.

3. Canopy: Once the leaves were "scabbed", we placed them in the scene, which started off as an artificial image of soil, generated in the same way as scab texture. The goal was to place them with certain occlusion and recovery patterns so that they matched real images as closely as possible. Our approach consisted in simulating apple "plants". We defined a certain number of randomly placed seeds, and iteratively added around each of these seeds a leaf. The placement of a leaf on a given seed was inspired by biological knowledge and observation of real images. In particular, angles between leaves of a plant are well defined because it is known that a plant tries to expose as many as possible to sunlight. This procedure resembled the "dead leaves model", a process proved to simulate images with statistical properties similar to the ones in natural images, such as scale invariance of statistics (Ruderman and Bialek, 1994; Lee et al., 2001).

To further account for the third dimension of the scene, the top halves of leaves were bent by homography. Moreover, leaves placed later in the iterative process were made bigger and more luminous than former ones, in order to simulate perspective. Luminosity and size changes were handled in an affine way, meaning leaves grew linearly brighter and bigger as they were placed closer to the top. Leaf luminosity was changed by multiplying the intensity of images.

Results can be seen in Fig. 4. Note that the simulator also generated annotation (pixels that represent scab) in addition to the IR images.

Ablation study: To learn more on which parts of the simulation were actually helping on pre-training, we conducted an ablation study, "knocking out" specific features of the simulator described in Section 2.3.2. We then pre-trained the model with these "ablated" images. The list of functions we isolated and knocked out one by one, keeping the others constant, was the following:

- 244 • Structure: leaves were kept in layers, but not organized in "apple plants". They were instead randomly positioned
245 and rotated.
- 246 • Texture: scab texture was generated using only mean and standard deviation of real scab texture images and not
247 autocorrelation (1), or only autocorrelation without the first order statistics (2).
- 248 • Size: leaves from different layers were not resized as a function of the layer they were in.
- 249 • Lighting: leaves from different layers' gray level intensity were not changed as a function of the layer they were
250 in.
- 251 • Bend: leaves were not bent by homography.
- 252 • Gradient: Scab texture was applied to leaves without any "gradient" effect.
- 253 • Vignetting: no vignetting effect was added to the image.
- 254 • Leaf species : leaves were not apple tree ones (*Malus pumila*), instead **mulberry** (*Broussonettia papyrifera*) ones.

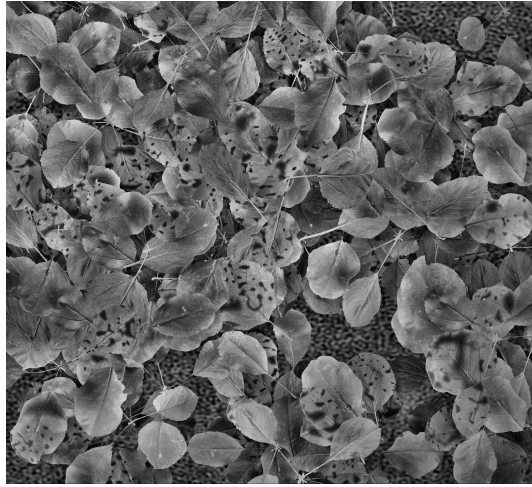


Figure 4: An example of an image coming from the model-based simulator.

255 2.3.3. GAN simulator

256 Another way to generate more images from a train set is by using GANs. We used Deep Convolutional GAN
257 (DCGAN) (Radford et al., 2015) architecture, with the Wasserstein GAN (WGAN) (Arjovsky et al., 2017) training
258 process. We tried implementation techniques to improve training stability (Salimans et al., 2016). The only change
259 that improved results was increasing the number of filters in the generator part of the network while leaving the
260 discriminator's number of filters to their default value. Training was done for approximately 100,000 generator
261 iterations, where the loss function did not decrease anymore and simulated images were stable. The difference with
262 traditional image generation by GAN is that the training images we used were concatenated images of IR images
263 (dimension $H \times W \times 1$) and their binary manual annotation (dimension $H \times W \times 1$), yielding an "annotated image"
264 (dimension $H \times W \times 2$) for the GAN to train on. Thus, the annotation of a given image was simply encoded as the
265 second channel of that image. There was therefore no need to manually create label images for generation afterward
266 (Zhu et al., 2018): a GAN trained in such a way is capable of generating IR images and their associated annotation.
267 The algorithm is illustrated in Fig. 5.

268 2.4. Evaluating the simulated images

269 Once simulated images were created using the simulators described in the previous section, we added them to
270 the train set. We studied two ways of using this additional data to help with the training process. To our knowledge,
271 there has been no research showing that one way of using additional data is more efficient than the other.

- 272 • Combining: Create a new train set by adding simulated data to the train set: **this** seemed like the most natural
273 way of augmenting the train set, as the network would train simultaneously on real and synthesized images.
274 (Ward et al., 2018). The training batches were constrained to contain 50% of real images and 50% of synthesized
275 images. We did not enforce such constraints and rely on the stochasticity of batch formation to generate well-
276 mixed batches.

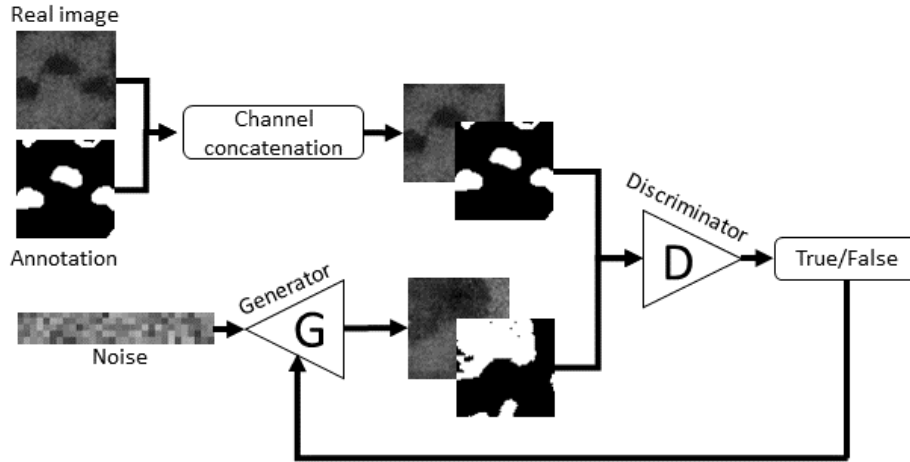


Figure 5: Our GAN-based algorithm description. G and D refer to generator and discriminator network, respectively. Grayscale images are IR images while binary images represent their corresponding annotation.

- **Pre-training:** Train a network on simulated data and fine-tune it on real images. Using features learned on a dataset to infer on another one, with or without fine-tuning, is one of the earliest practices of neural network training (Yosinski et al., 2014). For example, many applications use weights initialized to ones learned on the ImageNet dataset (Deng et al., 2009; Sharif Razavian et al., 2014). This is based on the idea that most natural images contain similar local patterns, and therefore neurons trained on a generic classification task such as classification on ImageNet are useful for a wide variety of applications. However, we hypothesized that one can achieve better results on a given segmentation task when pre-training on a task close to the final objective instead of on a more generic task.

For all "pre-training" experiments, we used datasets with 5,000 training images. For all "combining" experiences, we tried 4 different simulated dataset sizes corresponding to 0.25, 1, 4 and 8 times the number of real images in the dataset. Results presented in section 3 are shown for the quantities of images that yielded the best scores.

2.5. Comparing with real annotation

To compare the performance of our simulators with real data, we computed segmentation scores of the full dataset of 535 images. More precisely, we created from this full dataset 10 reduced train sets nested into each other, *i.e.* $d_1 \in d_2 \in \dots \in d_9 \in d_{full}$, where d_{full} is the dataset of 535 images. We then trained the network on each of the reduced train sets and computed the segmentation scores on the same test set. To complete our analysis, we computed the score on these different train set sizes for different augmentation strategies:

- case 1 : Only pre-training on ImageNet, with no other augmentation (base case);
- case 2 : Pre-training on ImageNet and using (online) SDA, which is still a basic use case of many segmentation routines;
- case 3 : Pre-training on data generated by the model-based simulator and using (online) SDA, which is the scenario we introduced in this paper that provided the highest performance.

3. Results

In this section, we present segmentation scores obtained after training the network presented in Section. 2.2 on the dataset presented in Section. 2, using the various data simulation strategies described in Section. 2.3. Results for a given strategy vary from a run to another, because of the stochasticity of neural network initialization, accentuated by the small dataset sizes. Therefore, all results are averaged over 30 runs, such that score variation between means of batches of 30 runs was under 0.5%. The results are presented in Table 1. With no augmentation whatsoever, the score was 0.424.

Table 1: Segmentation scores for the different dataset augmentation strategies. The ImageNet dataset is typically used for pre-training, and there was no sense in combining it with our dataset as images classification task wildly differ from our case.

Helper dataset	Combining	Pre-training
None	0.424 ± 0.013	
ImageNet	-	0.472±0.009
SDA (offline)	0.559±0.007	0.519±0.006
SDA (online)	0.574±0.010	-
Model-based simulator	0.509±0.005	0.602±0.007
GAN-based simulator	0.496±0.011	0.494±0.010

307 Results for the ablation study of the model-based simulator are presented in Fig. 6.

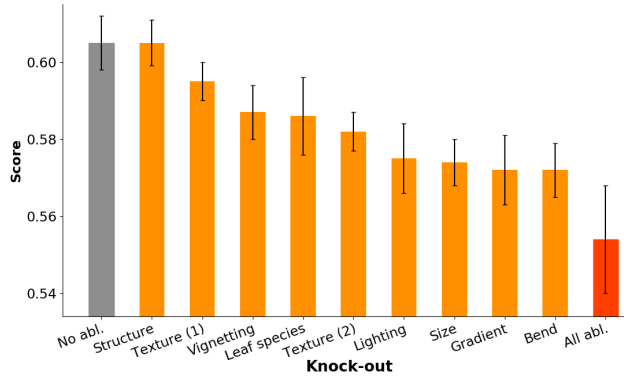


Figure 6: Segmentation scores for the different configurations of the ablation study of the model-based simulator. The X-axis indicates the specific feature of the simulator that was "knocked out" for that experiment. "No abl." refers to an experiment with no ablations (using the full simulator, as in Table 1). "All abl." refers to an experiment where all ablations described in Section 4.2 were applied at once.

308 We also tested combinations of augmentation strategies: namely adding online data augmentation to train sets
 309 partly composed of model-based or GAN-based data.

Table 2: Segmentation scores for the different dataset augmentation strategies combined.

Helper dataset	Combining	Pre-training
Model-based simulator + SDA (online)	0.580±0.004	0.643±0.005
GAN-based simulator + SDA (online)	0.599±0.008	0.529±0.007

310 Results of the comparison of data simulation and actual real data annotation are presented in Fig. 7.

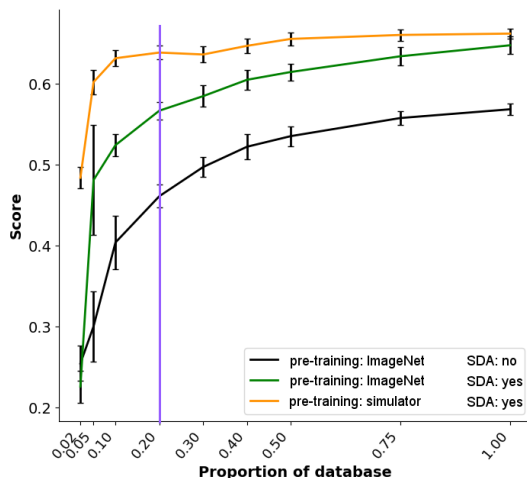


Figure 7: Impact of train set size on segmentation score. The X-axis indicates the proportion of train set kept for training from the full train set (535 images), while validation and test set stay the same size. The Y-axis indicated the segmentation score on the test set. The sub-train set presented in this article up to this point used 0.20 of the train set (marked by the vertical purple line). The different curves correspond to different types of augmentations, as indicated in the legend.

4. Discussion

Let us call the strategy of pre-training on ImageNet the "base case", as it is a widely used technique in machine learning algorithms. The result for the base case was a score of 0.472 (Table 1). This value corresponded to typical values of precision and recall of about 0.533 and 0.424, respectively. For all of our results, we observed relative values of precision and recall close to this case (values are quite close, with precision being higher than recall). To avoid overloading the presentation of the results, precision and recall are not mentioned in Tables 1 and 2. Segmentation suffered from both low precision and low recall. By observing outputs from the model, we noted that low recall came from entire scab lesions being missed. Low precision came from a segmentation "halo" of false positives around detected regions. A possible explanation for this effect was class balancing (see section 2.2), which strongly penalized false negatives to the point where the network became "overcautious" and classified all areas where the decision was difficult (in particular pixels in the aforementioned "gradient zone" (section 2.1)) as "scab", leading to these false positives "halos". Let us now discuss the impact of the different simulation strategies, compared to this base case. While comparing results, it is worth noting that experiments ran with 30 runs led to a standard deviation of the score typically under 1% (Table 1).

4.1. Standard data augmentation

SDA yielded very convincing progress for this dataset compared to the base case, especially in configurations where it was used directly in combination with our original images. In the offline version, the improvement compared to the base case was 8%. We suspect that with more images, we could have reached similar scores to the online version, which yielded a 10% improvement. This confirmed our intuition that SDA was well adapted to our case: a training image flipped in any direction, rotated in any direction, blurred and/or with changed perspective led to an image that could very much resemble images from the test set. The important positive impact of SDA on plant dataset segmentation is consistent with previous findings (Pawara et al., 2017).

4.2. Model-based simulator

Using the model-based simulator yielded interesting results: the simulated images were efficient in a pre-training scenario (+12%). This showed that images outputted by the simulator were at least somewhat realistic from the networks point of view. However, results when using these simulated images in a combined train set were not much higher than the ones obtained on the base case (+3%). The score difference of using these model-based simulated images as pre-training or combining was interesting, as it illustrated the important difference between these two methods of using extra data. In our case, model-based simulated images were very useful in pre-training, meaning they were efficient in leading training of weights "on the right tracks", but were too different from real data to be beneficial if added directly in the train set.

The ablation study (Fig. 6) gives us some insight on the usefulness of our simulator, by showing at which point each of the simulator's features was useful. The "gradient" effect of scab texture and leaf homography seemed to

345 have the most impact (3% difference with the full simulator) while placing leaves in "plant" structure seems to have
346 been almost useless. It is worth noting that even with all the features tested in this ablation study taken out at once
347 ("All" result in Fig. 6), results were still an 8% gain from the base case. It seems that the mere fact of pre-training
348 on images of real leaves greatly helped the network when training on the scab dataset. We hypothesized that the
349 success of this "simple" pre-training could be explained by the same reasons pre-training on ImageNet helps most
350 training tasks: **features** learned when pre-training on leaves seem generic enough to be useful to any "plant-related"
351 training tasks.

352 *4.3. GAN-based simulator*

353 Compared to the above augmentations, improvement by the GAN generator was sizably smaller: 2% regardless
354 the way data is used. One possible cause for the smallness of this effect was the difficulty for the GAN to converge
355 to convincing images. Fig. 8 presents different experiments on training the GAN on the train set, which can be seen
356 as intermediary steps to our full training process. While the visual quality of a simulated image did not necessarily
357 strictly correlate with the image's capacity to improve the segmentation process of a real dataset, we believed that
358 the former was a proxy for the latter.

359

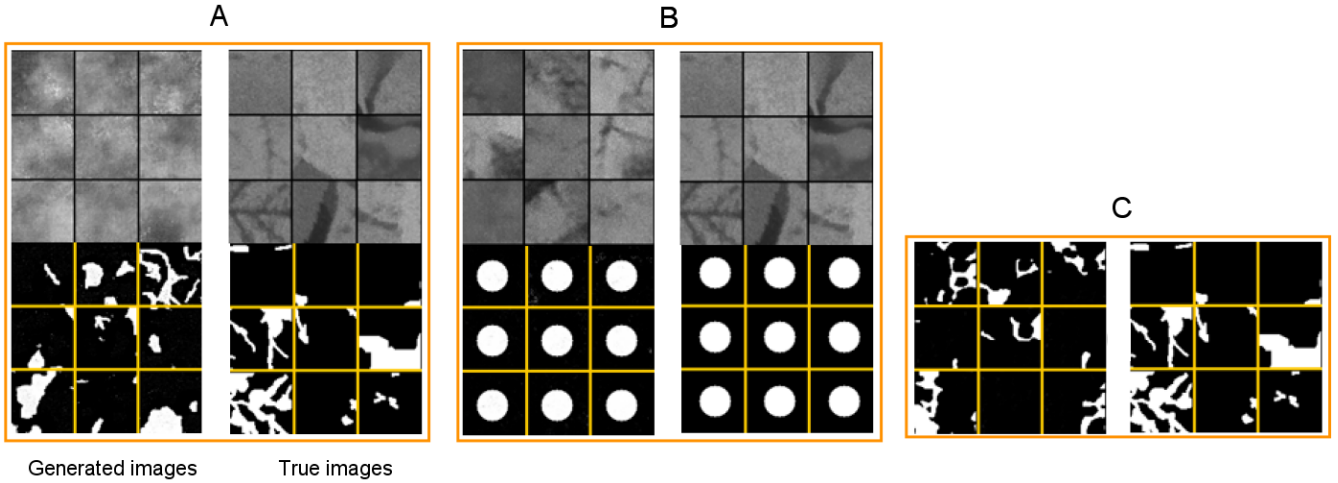


Figure 8: Different experiments (orange squares) to study the difficulties of the GAN based simulator. Each experiment shows 9 examples of original images the GAN trained on (right) and 9 examples of the generator’s output after convergence (left). When two blocks of 9 images are presented on top of another (exps. **A** and **B**), they represent an annotated (2-channel) image. Each binary image corresponds to the annotation of the IR image at a similar position. Exp. **A** is a case we present in this article: training on the annotated image. In Exp. **B**, we replaced annotations by “dummy” binary channels, all representing a white circle on black background. Exp. **C** presents GAN training done only on annotations.

360 The generation presented in this article is shown in Fig. 8 **A**: simulated images seemed relevant because they
 361 captured global distribution and structure of IR images and created annotations that seemed realistic when consid-
 362 ered on their own. However, they were not satisfactory in the sense that they failed to represent fine structures of
 363 original IR images (leaf veins and contours), they had artifacts and annotation structures seemed partly unrelated
 364 to their associated IR image. Fig. 8 **B** shows an experiment where the IR images were associated with a “fake”
 365 annotation unrelated to the images. Results were visually satisfying, indicating that our GAN algorithm was capa-
 366 ble of generating 2 channel images including a grayscale one and a binary one. Fig. 8 **C** presents GAN-simulated
 367 images from the binary images only. The GAN seemed to be perfectly capable of generating such binary structures.
 368 Thus, the difficulty resided in generating an IR image and a *relevant* binary channel, *i.e.* an annotation, as we
 369 saw with the difference between Fig. 8 **A** and **B**. Our intuition was that in order to create these annotated images,
 370 the GAN’s discriminator’s role was much more complex than when creating “simple” grayscale images. When the
 371 discriminator had to decide whether an annotated IR image came from a real or fake distribution, it had to judge
 372 (among other things) if the two channels were consistent with each other, *i.e.* it considered the image as a whole. In
 373 order to be able to decide about that, the discriminator had to find features in the IR channel that could robustly
 374 explain pixel values in the annotation channel: in other words, the discriminator had to act as a segmentation
 375 network. We believed that this is what made the task difficult for the GAN system and explained poorer results in
 376 exp. **A**.

377 4.4. Multiple augmentations

378 Our best results were obtained when combining pre-training on the model-based simulator and adding online
 379 SDA (Table 2): a 17% increase compared to the base case. This showed that different methods of augmentation
 380 could be combined very efficiently.

381 4.5. Comparing with real annotation

382 Fig. 7 shows the network’s scores on train datasets of varying size, with or without the augmentation techniques
 383 discussed in this article. We made the following observations. Firstly, the gains from the different augmentation
 384 strategies we observed in this article (vertical purple line on Fig. 7) also appeared for other train set sizes. Next,
 385 regardless of the chosen augmentation strategy, segmentation score as a function of train set size was similar to a
 386 logarithmic growth. This was quite intuitive : when showing new images to a network learning on a very small
 387 train set, these new images were very likely to show new patterns and structures, never seen before by the network.
 388 The more images were added, the more the network had seen “everything there is to see”, and the less it learned
 389 from new images. A corollary of this kind of growth was that scores tended to saturate. With the full train set,
 390 we seemed to have not reached full saturation yet. Interestingly, the gap between cases 1 and 2 did not seem to
 391 decrease as new images were added. This was a surprising result for us, as we would have expected the curves
 392 to converge as we added new images (SDA becoming less and less useful as we cover more and more cases), as it

was the case for the gap between cases 2 and 3. It seemed here that SDA may enable to get on the whole better results compared to adding some new "real" data. Especially for small train set sizes, we saw the usefulness of data augmentation strategies for saving annotation effort. For example, to get a score similar to the one obtained by annotating 150 images (30%) of the train set in case 1, one would need only about 30 images in case 2 and 10 in case 3. Finally, regardless of the type of augmentation, we found these curves to be interesting from a theoretical and practical point of view. Consider the case where one computes one of these curves for a dataset A , and also wishes to work on a dataset B that seems quite similar to A according to certain metrics (in our example, if A is our scab dataset, a similar B dataset could be images of another type of biotic stress on another plant species' leaves for example, with camera position and plant density close to A). Then one could assume that the "data utility curve" of A and B are similar, and therefore infer beforehand the sufficient quantity of B data to annotate for "best possible" segmentation performance. In the case of segmentation annotation, especially in a domain like agronomy where expert intervention is needed, knowing this kind of information can be very valuable.

4.6. Extension to other plant diseases

Many plant diseases show general structure close to the one of apple scab. For example, downy mildew of grapevines, late blight of potatoes and early symptoms of fusarium wilt will resemble small textured lesions appearing on leaves. The GAN-based simulator is in theory an image-agnostic tool, applicable to other segmentation tasks. Since GANs rely on finding some structure in data in a latent space, we believe the results presented here would be reproducible for these other diseases where image structure is close to our scab database. We believe that our model-based simulator could also be extended to these other types of plant diseases characterized by blobs on the leaves. Indeed, contrast, size, spatial distribution and texture of the lesions, size, density, and the species of leaves are all parameters of this simulator which could be tuned.

An extension of our simulator would be to be able to handle color images. While IR imagery was proved to be an appropriate modality for apple scab, standard RGB images can also provide valuable information on a wide range of plant diseases (Mahlein, 2016). In such RGB images, leaf edges are more clearly defined, perspective and leaf veins appear much more distinctly. An interesting perspective would be to provide both RGB and IR images to the network. Such multi-modal images are moreover becoming increasingly easier to acquire thanks to the development of convenient multi-modal image sensors for the plant science domain (Lowe et al., 2017).

5. Conclusion and future works

We have tackled the difficult and open problem in plant science of apple scab segmentation at canopy level. Like many segmentation tasks, the train set was quite small, and with this in mind we designed novel ways of generating new annotated data from the existing train set. This new data helped us achieve a much better segmentation, especially when different types of augmentation were combined.

Acknowledgements

Authors thank the PHENOTIC platform of INRA Angers, member of the PHENOME-EMPHASIS plant phenotyping network together with Charles-Eric Durel and Caroline Denance from INRA Angers for assistance in image acquisition of apple trees with scab in controlled environment. Clément Douarre gratefully acknowledges ANRT for CIFRE PhD funding under 2017/0639.

Arjovsky, M., Chintala, S., Bottou, L., 2017. Wasserstein generative adversarial networks. In: International Conference on Machine Learning. pp. 214–223.

Badrinarayanan, V., Kendall, A., Cipolla, R., 2017. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence* 39 (12), 2481–2495.

Barth, R., IJsselmuiden, J., Hemming, J., Van Henten, E. J., 2018. Data synthesis methods for semantic segmentation in agriculture: A capsicum annum dataset. *Computers and Electronics in Agriculture* 144, 284–296.

Belin, É., Rousseau, D., Boureau, T., Caffier, V., 2013. Thermography versus chlorophyll fluorescence imaging for detection and quantification of apple scab. *Computers and Electronics in Agriculture* 90, 159–163.

Benoit, L., Benoit, R., Belin, É., Vadaine, R., Demilly, D., Chapeau-Blondeau, F., Rousseau, D., 2016. On the value of the kullback–leibler divergence for cost-effective spectral imaging of plants by optimal selection of wavebands. *Machine Vision and Applications* 27 (5), 625–635.

- 442 Benoit, L., Rousseau, D., Belin, É., Demilly, D., Chapeau-Blondeau, F., 2014. Simulation of image acquisition
443 in machine vision dedicated to seedling elongation to validate image processing root segmentation algorithms.
444 *Computers and electronics in agriculture* 104, 84–92.
- 445 Bowen, J. K., Mesarich, C. H., Bus, V. G., Beresford, R. M., Plummer, K. M., Templeton, M. D., 2011. *Venturia*
446 *inaequalis*: the causal agent of apple scab. *Molecular Plant Pathology* 12 (2), 105–122.
- 447 Brostow, G. J., Fauqueur, J., Cipolla, R., 2009. Semantic object classes in video: A high-definition ground truth
448 database. *Pattern Recognition Letters* 30 (2), 88–97.
- 449 Brownlee, J., 2015. Tactics to combat class imbalance in your machine learning dataset. <https://machinelearningmastery.com/tactics-to-combat-imbalanced-classes-in-your-machine-learning-dataset/>,
450 accessed: 2018-10-14.
451
- 452 Chéné, Y., Rousseau, D., Lucidarme, P., Bertheloot, J., Caffier, V., Morel, P., Belin, É., Chapeau-Blondeau, F.,
453 2012. On the use of depth camera for 3d phenotyping of entire plants. *Computers and Electronics in Agriculture*
454 82, 122–127.
- 455 Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B.,
456 2018. Semantic understanding of urban street scenes. <https://www.cityscapes-dataset.com/benchmarks/>,
457 accessed: 2018-10-12.
- 458 Delalieux, S., Somers, B., Verstraeten, W., Van Aardt, J., Keulemans, W., Coppin, P., 2009. Hyperspectral indices
459 to diagnose leaf biotic stress of apple plants, considering leaf phenology. *International Journal of Remote Sensing*
460 30 (8), 1887–1912.
- 461 Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., Fei-Fei, L., 2009. Imagenet: A large-scale hierarchical image
462 database. In: *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on. IEEE*, pp.
463 248–255.
- 464 Di Cicco, M., Potena, C., Grisetti, G., Pretto, A., 2017. Automatic model based dataset generation for fast and
465 accurate crop and weeds detection. In: *2017 IEEE/RSJ International Conference on Intelligent Robots and*
466 *Systems (IROS). IEEE*, pp. 5188–5195.
- 467 Frid-Adar, M., Klang, E., Amitai, M., Goldberger, J., Greenspan, H., 2018. Synthetic data augmentation using
468 gan for improved liver lesion classification. In: *Biomedical Imaging (ISBI 2018), 2018 IEEE 15th International*
469 *Symposium on. IEEE*, pp. 289–293.
- 470 Giuffrida, M. V., Scharr, H., Tsaftaris, S. A., 2017. Arigan: Synthetic arabidopsis plants using generative adversarial
471 network. In: *Proceedings of the 2017 IEEE International Conference on Computer Vision Workshop (ICCVW),*
472 *Venice, Italy. pp. 22–29.*
- 473 Goodfellow, I., Bengio, Y., Courville, A., Bengio, Y., 2016. *Deep learning. Vol. 1. MIT press Cambridge.*
- 474 Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y., 2014.
475 *Generative adversarial nets. In: Advances in neural information processing systems. pp. 2672–2680.*
- 476 Isola, P., Zhu, J.-Y., Zhou, T., Efros, A. A., 2017. Image-to-image translation with conditional adversarial networks.
477 In: *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). IEEE*, pp. 5967–5976.
- 478 Kamilaris, A., Prenafeta-Boldú, F. X., 2018. Deep learning in agriculture: A survey. *Computers and Electronics in*
479 *Agriculture* 147, 70–90.
- 480 Kumar, N., Belhumeur, P. N., Biswas, A., Jacobs, D. W., Kress, W. J., Lopez, I., Soares, J. V. B., October 2012.
481 *Leafsnap: A computer vision system for automatic plant species identification. In: The 12th European Conference*
482 *on Computer Vision (ECCV).*
- 483 Lee, A. B., Mumford, D., Huang, J., 2001. Occlusion models for natural images: A statistical study of a scale-
484 invariant dead leaves model. *International Journal of Computer Vision* 41 (1-2), 35–59.
- 485 Lowe, A., Harrison, N., French, A. P., 2017. Hyperspectral image analysis techniques for the detection and classifi-
486 cation of the early onset of plant disease and stress. *Plant methods* 13 (1), 80.

- 487 Mahlein, A.-K., 2016. Plant disease detection by imaging sensors—parallels and specific demands for precision
488 agriculture and plant phenotyping. *Plant Disease* 100 (2), 241–251.
- 489 Matheron, G., Matheron, G., Matheron, G., Matheron, G., 1975. Random sets and integral geometry.
- 490 McCombs, S., 2005. Intro to procedural textures. [http://www.upvector.com/?section=Tutorials&subsection=](http://www.upvector.com/?section=Tutorials&subsection=IntrotoProceduralTextures)
491 [IntrotoProceduralTextures](http://www.upvector.com/?section=Tutorials&subsection=IntrotoProceduralTextures), accessed: 2018-10-10.
- 492 Minervini, M., Scharr, H., Tsafaris, S. A., 2015. Image analysis: the new bottleneck in plant phenotyping [appli-
493 cations corner]. *IEEE signal processing magazine* 32 (4), 126–131.
- 494 Mirza, M., Osindero, S., 2014. Conditional generative adversarial nets. arXiv preprint arXiv:1411.1784.
- 495 Neff, T., Payer, C., Štern, D., Urschler, M., 2017. Generative adversarial network based synthesis for supervised
496 medical image segmentation.
- 497 Noh, H., Hong, S., Han, B., 2015. Learning deconvolution network for semantic segmentation. In: Proceedings of
498 the IEEE international conference on computer vision. pp. 1520–1528.
- 499 Oerke, E.-C., Fröhling, P., Steiner, U., 2011. Thermographic assessment of scab disease on apple leaves. *Precision*
500 *Agriculture* 12 (5), 699–715.
- 501 Pawara, P., Okafor, E., Schomaker, L., Wiering, M., 2017. Data augmentation for plant classification. In: Interna-
502 tional Conference on Advanced Concepts for Intelligent Vision Systems. Springer, pp. 615–626.
- 503 Pound, M. P., Atkinson, J. A., Townsend, A. J., Wilson, M. H., Griffiths, M., Jackson, A. S., Bulat, A., Tzimiropou-
504 los, G., Wells, D. M., Murchie, E. H., et al., 2017. Deep machine learning provides state-of-the-art performance
505 in image-based plant phenotyping. *GigaScience*.
- 506 Pradal, C., Dufour-Kowalski, S., Boudon, F., Fournier, C., Godin, C., 2008. Openalea: a visual programming and
507 component-based software platform for plant modelling. *Functional plant biology* 35 (10), 751–760.
- 508 Radford, A., Metz, L., Chintala, S., 2015. Unsupervised representation learning with deep convolutional generative
509 adversarial networks. arXiv preprint arXiv:1511.06434.
- 510 Ruderman, D. L., Bialek, W., 1994. Statistics of natural images: Scaling in the woods. In: *Advances in neural*
511 *information processing systems*. pp. 551–558.
- 512 Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., Chen, X., 2016. Improved techniques for
513 training gans. In: *Advances in Neural Information Processing Systems*. pp. 2234–2242.
- 514 Sharif Razavian, A., Azizpour, H., Sullivan, J., Carlsson, S., 2014. Cnn features off-the-shelf: an astounding baseline
515 for recognition. In: *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*.
516 pp. 806–813.
- 517 Ubbens, J., Cieslak, M., Prusinkiewicz, P., Stavness, I., 2018. The use of plant models in deep learning: an
518 application to leaf counting in rosette plants. *Plant methods* 14 (1), 6.
- 519 Ward, D., Moghadam, P., Hudson, N., 2018. Deep leaf segmentation using synthetic data. arXiv preprint
520 arXiv:1807.10931.
- 521 Yosinski, J., Clune, J., Bengio, Y., Lipson, H., 2014. How transferable are features in deep neural networks? In:
522 *Advances in neural information processing systems*. pp. 3320–3328.
- 523 Zhu, Y., Aoun, M., Krijn, M., Vanschoren, J., Campus, H. T., 2018. Data augmentation using conditional generative
524 adversarial networks for leaf counting in arabidopsis plants. *Computer Vision Problems in Plant Phenotyping*
525 (CVPPP2018).