

# On Domain-Specific Pre-Training for Effective Semantic Perception in Agricultural Robotics

Gianmarco Roggiolani

Federico Magistri

Tiziano Guadagnino

Jan Weyler

Giorgio Grisetti

Cyrill Stachniss

Jens Behley

**Abstract**—Agricultural robots have the prospect to enable more efficient and sustainable agricultural production of food, feed, and fiber. Perception of crops and weeds is a central component of agricultural robots that aim to monitor fields and assess the plants as well as their growth stage in an automatic manner. Semantic perception mostly relies on deep learning using supervised approaches, which require time and qualified workers to label fairly large amounts of data. In this paper, we look into the problem of reducing the amount of labels without compromising the final segmentation performance. For robots operating in the field, pre-training networks in a supervised way is already a popular method to reduce the number of required labeled images. We investigate the possibility of pre-training in a self-supervised fashion using data from the target domain. To better exploit this data, we propose a set of domain-specific augmentation strategies. We evaluate our pre-training on semantic segmentation and leaf instance segmentation, two important tasks in our domain. The experimental results suggest that pre-training with domain-specific data paired with our data augmentation strategy leads to superior performance compared to commonly used pre-trainings. Furthermore, the pre-trained networks obtain similar performance to the fully supervised with less labeled data.

## I. INTRODUCTION

Sustainable crop production is fundamental to meet the increasing request for food, fuel, and fiber. It, however, must become more effective to fulfill all demands. Furthermore, the lack of workers is a key challenge, which is even increased during the recent COVID pandemic. Robots are a crucial component to analyze and monitor plants in an automated way [28], followed by targeted fertilizing and/or protection [13]. Before targeted actions can be performed, the underlying perception problems need to be solved. Deep learning approaches improved the performance of these systems using large neural networks. Such networks for object detection or semantic segmentation can help the farmers to evaluate the status of the plants [17][27], spot and locate weeds [39], detect plant diseases [15], and understand the growing conditions among different areas of the field [40].

G. Roggiolani, F. Magistri, T. Guadagnino, J. Weyler, J. Behley, and C. Stachniss are with the University of Bonn, Germany. C. Stachniss is additionally with the Department of Engineering Science at the University of Oxford, UK, and with the Lamarr Institute for Machine Learning and Artificial Intelligence, Germany. G. Grisetti is with La Sapienza University of Rome, Italy.

This work has partially been funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany's Excellence Strategy, EXC-2070 – 390732324 – PhenoRob, and by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under STA 1051/5-1 within the FOR 5351 (AID4Crops).

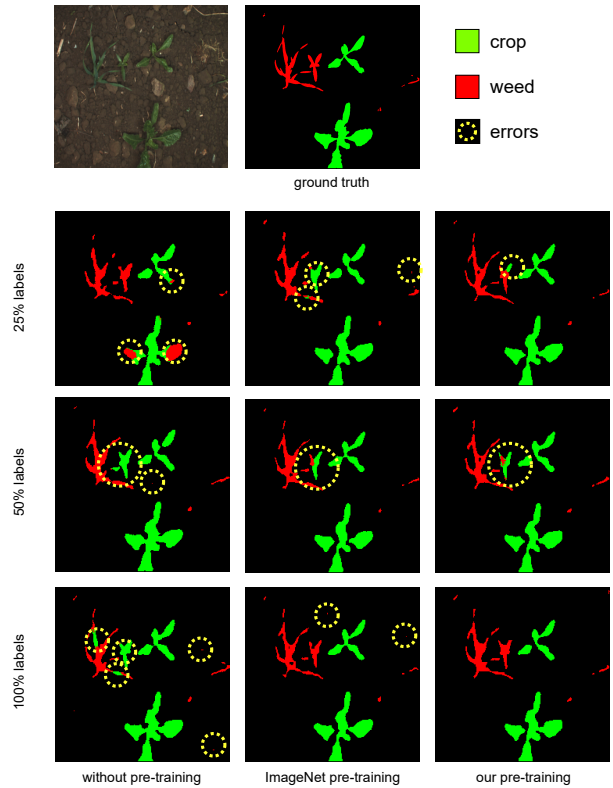


Fig. 1: Results on semantic segmentation with different amounts of data and pre-trainings. We achieve better or comparable results with  $\frac{1}{4}$  of the epochs and  $\frac{1}{1000}$  of the images. The dotted circles highlight the errors in the results.

Deep learning-based approaches, however, generally need a large amount of labeled data, which is hard to get because it needs time and specialized workers. Some researchers use semi-supervised approaches to reduce the need for labeled images [26] or leverage background knowledge [30]. Recent work in self-supervised pre-training showed promising results, where we can pre-train a network without relying on supervision by labels. Pre-training on the ImageNet dataset [9] is a common way to reduce the number of training samples and the training time for the network to converge. Several other methods use a contrastive loss and strong augmentation techniques [4][6][16][19][46]. Embeddings produced in such a way are more robust to the augmentations applied. However, the applied data augmentations need to be selected carefully so as not to lose relevant features. In this work, we study self-supervised representation learning to improve the perception of agricultural robots.

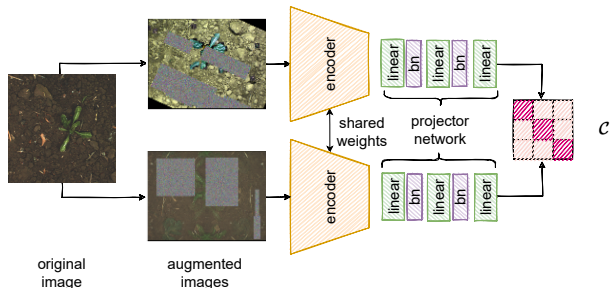


Fig. 2: For each image we build two augmentations to produce two embeddings. We fill a matrix  $C$  with their cross-correlation and then compute the loss, that forces  $C$  to be an identity.

The main contribution of this paper is providing a pre-training strategy for the plant domain, which will reduce the number of labeled images needed and a newly defined augmentation policy. We study existing augmentation and propose domain-specific ones to boost performance. Specifically, we target semantic and leaf instance segmentation and investigate how self-supervised pre-training on domain-specific data leads to better models that can learn with less labeled data. In sum, we make three key claims: (i) pre-training on datasets of the same domain improves the performance of the downstream tasks, (ii) using domain-specific pre-training can further reduce the number of labeled images needed, and (iii) the augmentation policy needs to be domain-specific and take into account the order and design of the augmentations. Our code is available at <https://github.com/PRBonn/agri-pretraining>.

## II. RELATED WORK

Robotic applications in agriculture aim at improving field monitoring and interventions diminishing the use of agricultural chemical inputs and production costs [32] [38]. Semantic segmentation and instance segmentation are two key steps for weed control and plant phenotyping. The majority of recent solutions for these tasks employ large convolutional neural networks [23] [24] [32], requiring large amounts of labeled data. Some of them exploit spatial information about the fields [25] or vegetation indexes [30], while others focused on the architecture side, as Potena et al. [31] where they use two CNNs to detect the vegetation and then classify it, Weyler et al. [39] which employs a Feature Pyramid Network to detect the instances, or in the works from Buzzy et al. [3] and You et al. [45] that use deep neural networks.

Two main paths have been investigated to reduce the number of labeled images: domain adaptation and pre-training. In the domain adaptation setting, the network learns how to perform a task on a source domain and is expected to have a good performance on a different one. Approaches often use generative adversarial networks [22] [44]. Their application on the plant domain shows already promising results for classification [14], object counting [1], and object detection [18]. However, these methods require training the network on a source domain dataset for which we might need vast amounts of labeled images.

In contrast, pre-training aims at initializing the weights of a neural network, such that it needs fewer labels and converges faster. Erhan et al. [10] shows that unsupervised pre-training guides the networks towards the minimum, from where supervised training can proceed faster and with less available data. In literature, pre-training on ImageNet is a common choice [35], since it is big enough to provide a good initialization for various tasks and domains.

Lately, self-supervised pre-training received increasing attention due to its promising results on several downstream tasks compared to supervised pre-training [6] [19] [7] [46]. Especially, contrastive approaches firstly used large memory banks [43] that were later replaced by a momentum encoder [19] to reduce the required memory to store negative samples. In particular, Chen et al. [6] introduced a projection head for learning embeddings of positive and negative examples that showed superior performance and were later integrated into momentum contrast [7]. They also investigate various data augmentation techniques and show the relevance of different augmentations, but also an order dependence of the augmentations. Grill et al. [16] build upon earlier work [6] and proved that negative examples are not strictly necessary for pre-training. Zbontar et al. [46] extends this idea of using only positive examples by measuring the cross-correlation between two augmented views of the same image.

Recently, He et al. [20] challenged the need for pre-training and showed that, given enough iteration and data, a randomly initialized network can reach the same performance. Their work also confirms that pre-training is an effective way to reduce the need for labeled data and time to converge. McCormac et al. [29] compare the model pre-trained on ImageNet versus the model pre-trained on their synthetic RGB-D dataset, whose domain is aligned with the target dataset and task. Their domain-specific pre-training performs better than the ImageNet pre-training, especially when using depth information.

In contrast to the supervised approaches, we aim at exploiting the large amount of unlabeled images that we can record using robotic systems. For the self-supervised pre-training, we use Barlow Twins [46]. We evaluate the importance of pre-training directly in the agricultural domain, instead of adopting a general pre-training on ImageNet. In our work, we show the advantages of domain-specific augmentations, for which there is not much literature [8] [34] and we examine how their order and application may influence the performance of the final system on the downstream task.

## III. OUR APPROACH

We aim at learning an abstract representation that will serve as a starting point for further learning tasks in the domain of the perception of plants. By deploying robots in the fields, we can quite easily collect a large amount of *unlabeled* data. This offers the potential to build systems to train networks in a self-supervised fashion. For our perception task, we pre-train the network encoder following Barlow Twins (BT) proposed by Zbontar et al. [46]. We decided for BT to avoid the problem of negative pairs since in



Fig. 3: We applied all of our augmentations to a single image to show one possible outcome for each one. In our pre-training strategy they are applied sequentially producing even more variations of the input.

the agricultural domain most of the images represent plants; the pre-training approach, as well as the architecture for the tasks, is not our main focus. We propose a domain-specific augmentation policy to boost the performance of the final system. We evaluate our pre-training on semantic and leaf instance segmentation.

### A. Barlow Twins

BT learns representations in a self-supervised fashion via redundancy reduction. Here, we briefly summarize its relevant parts and refer for more details to the original paper [46]. It uses a siamese network with shared weights, which can be seen in Fig. 2. The two inputs are two different augmentations of the same input image. The encoder is a ResNet50 [21] without the final classification layer, followed by a projector network. The projector network has two identical blocks — linear layer, batch normalization, and rectified linear units — followed by one linear layer.

Zbontar et al. build for each input image  $I$  two augmented views  $I_1, I_2$  that are fed into the network to produce two distinct embeddings  $z_1, z_2 \in \mathbb{R}^D$ . They compute the loss directly on  $z_1$  and  $z_2$ . The first step is to construct the cross-correlation squared matrix  $\mathcal{C} \in \mathbb{R}^{D \times D}$  from the embeddings normalized over the batch dimension. This matrix has values between  $-1$  (anti-correlation) and  $1$  (correlation). Then, the loss is:

$$\mathcal{L}_{BT} \triangleq \sum_i (1 - \mathcal{C}_{ii})^2 + \lambda \sum_i \sum_{j \neq i} \mathcal{C}_{ij}^2, \quad (1)$$

where  $\lambda$  is a weight to trade-off two parts: the invariance term, which forces the diagonal elements of  $\mathcal{C}$  to be 1, and the redundancy reduction term, which forces all the non-diagonal elements of  $\mathcal{C}$  to be 0.

### B. Augmentations

Augmentations play a fundamental role in self-supervised learning. The stronger they are, the more the network focuses on relevant and stable features to represent the images. We use our augmentation policy as common in the literature [16] plus domain-specific knowledge. In Fig. 3, we show the result of each augmentation applied to one sample image for illustration purposes. Our augmentations are:

1) *Affine Transformation*: The affine transformation,  $\mathbf{T}_{\text{affine}} \in \mathbb{R}^{3 \times 3}$  rotates, translates, scales, and shears the input image. It makes the network invariant to such transformations which are common when working with robots of different sizes and cameras. More specifically,  $\mathbf{T}_{\text{affine}}$  is given by

$$\mathbf{T}_{\text{affine}} = \begin{bmatrix} \mathbf{A} & \underline{\mathbf{t}} \\ \mathbf{0}^T & 1 \end{bmatrix}, \quad (2)$$

where  $\mathbf{A} \in \mathbb{R}^{2 \times 2}$  contains an isotropic scaling factor in  $[0.5, 2]$ , a rotation in  $[-\pi, \pi]$  and the shearing along the two axes randomly sampled in  $[0.25, 0.75]$ , while  $\underline{\mathbf{t}} \in \mathbb{R}^2$  is a translation vector with each component  $t_x, t_y \in [-0.25 \cdot W, 0.25 \cdot H]$ , where  $W$  and  $H$  are the image width and height respectively.

2) *Color Jittering*: Color jittering changes the brightness, contrast, hue, and saturation of the image. Instead of the symmetrical range of values for the hue  $(-0.1, 0.1)$  from the literature, we use  $(0, 0.125)$  as range. This transformation is crucial for the classification of ill or damaged plants, where color is a dominant discriminator [37].

3) *Gaussian Blur*: We blur the image using a random standard deviation  $\in [0.1, 2]$ . The purpose of this augmentation is to help the network focus on the image structure across different scales and resolutions.

4) *Mixing*: Zhang et al. [47] propose to mix two images via linear interpolation, we instead use a single image  $I$ . We create two copies of the image  $I$ , one flipped on the x-axis  $I_x$  and one on the y-axis  $I_y$ . We sample each pixel in the augmented image from  $I, I_x$ , or  $I_y$  using a uniform probability over the three images. This can simulate motion due to the wind or water uptake, or holes eaten by insects.

5) *Random Erasing*: Random erasing [48] selects multiple rectangles inside of the image and substitutes the pixels' values with random values in  $[0, 255]$ . Given the minimum percentage of the image that has to be removed, it picks rectangles of different sizes and aspect ratios until the deleted area is at least the minimum desired area. We slightly change the implementation to enforce the use of multiple rectangles with respect to a big one. We use this augmentation to make the network less sensitive to occlusions and shadows.

6) *Background Invariance*: This transformation cuts plants from the current image and pastes them into a different soil background. Specifically, we perform the following steps:

(i) Compute a normalized image as

$$I_{\text{norm}}(u, v) = \frac{I(u, v) - \mu_I}{\sigma_I + \epsilon}, \quad (3)$$

where  $u, v$  are pixel coordinates,  $\mu_I \in \mathbb{R}^3$  and  $\sigma_I \in \mathbb{R}^3$  are mean and standard deviation of  $I$ , and  $\epsilon = 10^{-8}$ .

(ii) Compute the vegetation mask  $M$  following Woebbecke et al. [42]. Specifically,  $M$  is given by

$$M = 2I_G - I_R - I_B, \quad (4)$$

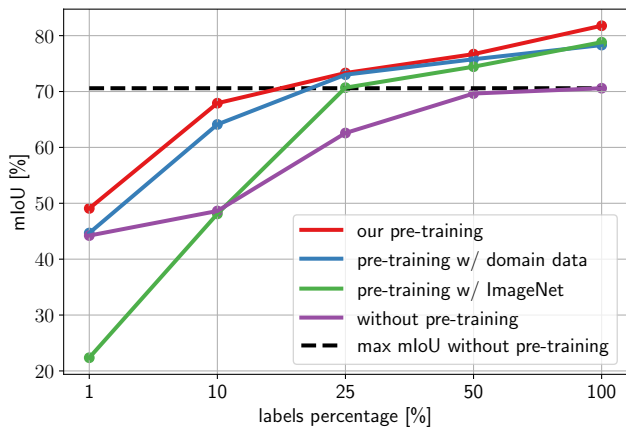


Fig. 4: Comparison of the mIoU with different amounts of labels after fine-tuning for 100 epochs (semantic segmentation). The number of images for each label percentage is: 14 for 1%, 140 for 10%, 362 for 25%, 724 for 50%, and 1450 for 100%.

where  $I_R, I_G$  and  $I_B$  are the color channels of  $I_{norm}$ .

- (iii) Convert  $M$  to a binary mask using a threshold  $\theta$ , i.e., all pixels above  $\theta$  are set to 1, the others are set to 0.
- (iv) Refine  $M$  using 2 rounds of erosion with kernel size  $(2, 2)$ , 4 rounds of dilation with kernel size  $(6, 6)$ .
- (v) Cut the vegetation and paste it at a random location on a random soil image from a dataset of images whose vegetation mask  $M$  is below a given threshold, i.e., less than 5% of the image.

All the augmentations are applied with a certain probability, tuned from the results in Sec. V-D.

#### IV. DOWNSTREAM TASKS

In an application, the pre-trained models are fine-tuned on specific downstream tasks. Since the pre-training approach is backbone invariant, we pre-trained ResNet50 for the semantic segmentation and an ERFNet [36] encoder for the leaf instance segmentation. For both tasks, we use ERFNet-like decoders.

##### A. Semantic Segmentation

Semantic segmentation predicts a class for each pixel of the image, in our application example, crop, weed, and background. The decoder outputs an image  $H \times W \times C_{out}$ , where  $C_{out}$  is the number of semantic classes,  $H$  and  $W$  are the height and width of the input image. Instead of connecting the decoder at the end of the ResNet50, we discard the last two layers to preserve more spatial information. We can initialize the remaining part with the pre-trained weights without changing anything. We follow Rahman et al. [33] to directly optimize the IoU.

##### B. Leaf Instance Segmentation

Leaf instance segmentation predicts a pixel-wise mask for each leaf. Such task allows discovering not only the shape and size of individual leaves but also counting them, which is fundamental in determining the growth stage of the plant [12]. We use the network and loss proposed by Weyler et al. [41]. One decoder predicts the center locations

TABLE I: Comparison of average precision (AP) and recall (AR) on plants (p) and leaves (l) for the three main approaches. The number of images for label percentage is: 7 for 1%, 74 for 10%, 186 for 25%, 373 for 50%, and 746 for 100%.

Pre-Training	$AP_p$ [%]	$AR_p$ [%]	$AP_l$ [%]	$AR_l$ [%]
100% of labels				
none	54.3	60.5	48.7	68.3
ImageNet	55.1	61.2	59.7	68.9
ours	<b>55.6</b>	<b>62.9</b>	<b>64.4</b>	<b>69.2</b>
50% of labels				
none	50.3	59.0	45.6	60.1
ImageNet	52.4	60.1	52.7	61.5
ours	<b>54.6</b>	<b>60.8</b>	<b>54.6</b>	<b>62.7</b>
25% of labels				
none	48.0	55.2	42.0	46.1
ImageNet	50.2	56.1	50.6	56.6
ours	<b>50.9</b>	<b>56.9</b>	<b>53.8</b>	<b>60.6</b>
10% of labels				
none	46.6	54.0	20.7	39.6
ImageNet	46.8	53.7	29.5	38.4
ours	<b>48.0</b>	<b>54.2</b>	<b>42.5</b>	<b>49.2</b>
1% of labels				
none	0.0	0.0	0.1	0.3
ImageNet	0.0	0.0	0.4	0.2
ours	<b>1.1</b>	<b>8.3</b>	<b>0.9</b>	<b>5.6</b>

of each leaf, the other the offsets pointing at the specific leaf and plant center plus clustering parameters for the post-processing. The predicted and the ground truth masks are fed to the Lovász Hinge Loss [2]. For more details, refer to the original paper [41].

#### V. EXPERIMENTAL EVALUATION

The focus of this work is showing that a domain-specific self-supervised pre-training together with augmentation policy has the potential to perform better than the commonly used supervised pre-training on ImageNet. We show this for images from the agricultural robotics domain. We present our experiments to show the capabilities of our pre-training and to support our key claims, which are: (i) pre-training on plant images improves the performance of the downstream tasks, (ii) using domain-specific pre-training can further reduce the number of labels needed, and (iii) the augmentation policy needs to be domain-specific and consider the design of each augmentation to apply them in the best order.

##### A. Experimental Setup

Our best model has the encoder pre-trained for 250 epochs, with batch size 128, learning rate  $2 \cdot 10^{-4}$ , and a weight decay of  $10^{-6}$ . We use 18,000 images from four different locations (Ancona in Italy, Bonn and Stuttgart in Germany, and Eschlikon in Switzerland, see also the dataset paper [5]) for pre-training. We compare it with a publicly available supervised pre-training on ImageNet. For the semantic segmentation task, we use a dataset that contains 2,148 images: 1,450 for training, 478 for validation, and 220 for testing. For the leaf instance segmentation, we use a dataset with 1,316 images; 746 for training, 292 for validation, and 278



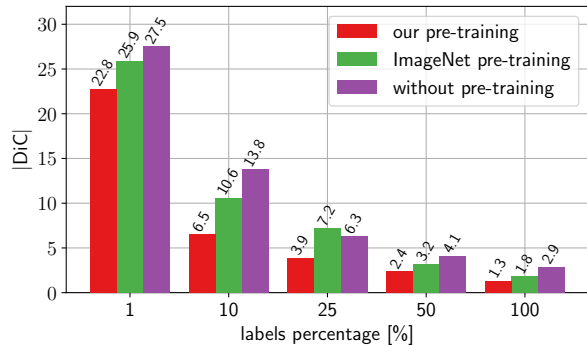


Fig. 5: The average  $|DiC|$  for the three approaches with increasing number of labeled images. The lower the better.

for testing. We pre-trained on a single NVIDIA RTX A6000 GPU and fine-tuned on a single Quadro RTX 5000 GPU. We plan to publish the datasets used both for pre-training and fine-tuning.

We evaluate the results on semantic segmentation using the mean intersection over union (mIoU) [11] over soil, crop, and weed. For the leaf instance segmentation task, we report the average precision (AP) and recall (AR), and the absolute difference in count ( $|DiC|$ ) of the leaves.

### B. Our Pre-training vs. Non-specific Pre-training

The first experiment analyzes how self-supervised domain-specific pre-training improves the effectiveness and decreases the labeled images, time, and computational resources needed. We fine-tuned our model and the model pre-trained on ImageNet with different amounts of labels.

**Semantic Segmentation.** Fig. 4 suggests that when using a sufficient number of labels, different pre-training strategies perform similarly. The fewer labels we use, the wider the gap. The ImageNet pre-training requires more labeled data to adjust to the agriculture domain. Our pre-training performs better requiring less data for pre-training i.e. 18,000 images against the 1,281,167 from ImageNet, and epochs i.e. 250 against 1,000. Only for this experiment, we also pre-trained on domain-specific data with the augmentations from the literature. The results confirm that domain-specific augmentations are a key component to obtain the best performance.

**Leaf Instance Segmentation.** Tab. I confirms the utility of pre-training on domain-specific data to boost the performance and reduce the number of labeled images needed. Our pre-training boosts every metric in every scenario. In Fig. 5, the difference in count shows a similar trend. When using less than 10 images none of the approaches can properly segment the leaves, but using 74 images (10%) our pre-training can already reduce the uncounted leaves to  $\approx 6$  per image (each image can have between 2 and 5 plants in it).

### C. Our Pre-training vs. No Pre-training

This experiment aims to compare the results and computational resources when using the pre-trained model with respect to training the network after a random initialization.

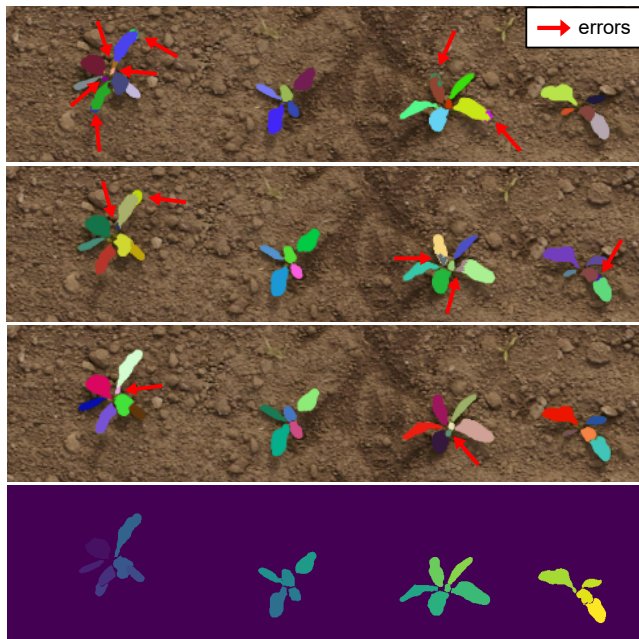


Fig. 6: Leaf instance segmentation with 50% of the labels. From top to bottom, we show the result with no pre-training, ImageNet and our pre-training, then the ground truth.

**Semantic Segmentation.** In Fig. 4, we see how pre-training boosts the performance when using the same routine and number of labeled images. Our pre-trained model performs better using up to 25% of the annotated data. With less than 10% of the labeled data, pre-training on ImageNet hurts the performance. The reason may be that the network expects to see objects from the ImageNet dataset distribution and requires the labeled data to adapt to the agricultural domain. Our pre-training does not suffer from this issue.

**Leaf Instance Segmentation.** Training the randomly initialized network with few labeled images deteriorates rapidly the performance. Fig. 5 and Tab. I show that we can obtain the same performance using half of the data. Our pre-training boosts all of the metrics, and it produces instance masks using only 7 images (1%).

### D. Relevance and Order of the Augmentations

We used a shorter training routine to analyze which augmentations work better for the plant domain. Each combination has been pre-trained for 50 epochs on a subset of the pre-training dataset, then fine-tuned on a subset of the dataset for semantic segmentation. Chen et al. [6] did a similar experiment whose results agree with ours.

In Fig. 7, we see that changing the transformation’s order impacts the mIoU and the training time. The combinations that took more time are those that make the task much harder, i.e., if we first apply color jittering and then background invariance, it will be challenging to correctly identify the plants, with the vegetation mask corrupted by the changes in color. Focusing on the highest-performing combinations we see that swapping them leads to lower values, confirming that the order in which they are applied is a key aspect when designing the augmentation policy. On the diagonal,

First Transformation	Second Transformation					
	Affine	Background Invariance	Color Jittering	Gaussian Blur	Mixing	Random Erasing
Affine	0.37	0.56*	0.41	0.36	0.42	0.43
Background Invariance	0.41	0.44	0.65	0.61	0.58	0.53**
Color Jittering	0.58	0.37**	0.48	0.51	0.42**	0.68*
Gaussian Blur	0.41	0.46	0.57	0.36	0.51**	0.38
Mixing	0.53*	0.38	0.63	0.55	0.41	0.7
Random Erasing	0.53	0.51	0.53***	0.45	0.35	0.46

Fig. 7: The mIoU for different combinations of transformations. We fine-tuned with minimum 100 epochs or until convergence: (\*) 10 extra epochs; (\*\*) 40 extra epochs; (\*\*\*) 100 extra epochs.

TABLE II: The mIoU [%] after fine-tuning (100 epochs on semantic segmentation) with 20, 40, 60, 80 and 100 epochs of pre-training using only the color transformation.

Approach	pre-training epochs				
	20	40	60	80	100
standard	23.44	23.77	23.89	19.92	13.89
our	17.61	21.11	31.28	32.52	42.80

where we apply only one augmentation, the mIoU values are all in the mid-lower range, the highest values being the strongest augmentations. This pattern is also visible in the other combinations; strong augmentations such as random erasing and mixing lead to better performance, not always at the price of longer training time.

As explained in Sec. III-B.2, we changed the usual parameters for the color jittering augmentation. To evaluate if this choice makes a difference we pre-trained our encoder only with color jittering and used the weights at different stages for the semantic segmentation task. We demonstrate in Tab. II that a longer pre-training period with the standard color augmentation degrades performance. One reason could be that the augmentation is so strong that the network does not take color into account anymore, making the semantic segmentation task harder. Our augmentation instead provides better performance the more the encoder is pre-trained.

### E. Influence of Pre-training Length

We pre-trained up to 500 epochs, evaluating intermediate checkpoints on the semantic segmentation task. We want to determine how many epochs are needed to achieve satisfying results and to find out how much we can improve by continuing training. This allows us to find a compromise between performance and computational resources. In Fig. 8 we show that the mIoU increases until 250 epochs, where

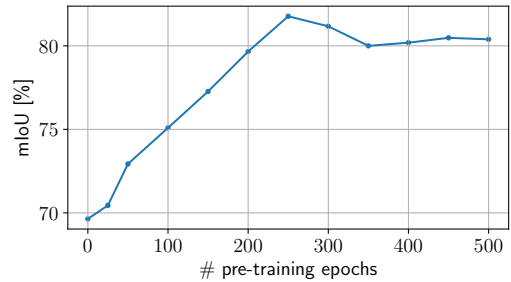


Fig. 8: The mIoU after fine-tuning for 100 epochs using different pre-training epochs. There is no improvement after 250 epochs. The starting value is obtained with a randomly initialized network.

TABLE III: Comparison of mIoU, mean precision (mP), and mean recall (mR) after fine-tuning for 100 epochs (semantic segmentation) with our augmentation’s probabilities vs. applying always all the augmentations.

Approach	mIoU [%]	mP [%]	mR [%]
all augmentations	78.09	88.61	87.68
our policy	<b>81.77</b>	<b>91.07</b>	<b>89.99</b>

we see a diminishing effect. Therefore, if not otherwise specified, we pre-train for 250 epochs in our experiments.

### F. Ablation on Augmentations’ Probabilities

Each augmentation is applied with a probability based on the results of the previous experiments to increase variance. For evaluation of the effectiveness of our policy, we compare it against a pre-training in which all augmentations are always applied.

The results in Tab. III show that the probabilities we assign to each augmentation result in higher performance on all metric evaluations. We propose to assign a probability of 1.0 to color jittering and random erasing, 0.9 to gaussian blur and mixing, and 0.8 to background invariance and affine transform.

## VI. CONCLUSION

In this paper, we presented an approach to exploit a vast quantity of unlabeled images from the agricultural domain to learn useful representations in a self-supervised fashion. Our experiments rely on domain-specific data and domain-specific augmentations during the pre-training. This allows us to successfully use our pre-training for different downstream tasks obtaining good performance using less labeled images. We implemented and evaluated our pre-trainings on two tasks, semantic and leaf instance segmentation in the agriculture domain. We compared our results with those obtained without pre-training and pre-training on ImageNet and supported all claims made in this paper. The experiments suggest that pre-training on a domain-specific dataset and exploiting domain knowledge to define the augmentation policy can reduce the number of labeled data required to achieve the same performances as without pre-training.

## REFERENCES

- [1] T. Ayalew, J. Ubbens, and I. Stavness. Unsupervised domain adaptation for plant organ counting. In *Proc. of the Europ. Conf. on Computer Vision (ECCV)*, 2020.
- [2] M. Berman, A.R. Triki, and M.B. Blaschko. The lovasz-softmax loss: A tractable surrogate for the optimization of the intersection-over-union measure in neural networks. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [3] M. Buzzy, V. Thesma, M. Davoodi, and J. Mohammadpour Velni. Real-time plant leaf counting using deep object detection networks. *Sensors*, 20(23):6896, 2020.
- [4] M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin. Unsupervised learning of visual features by contrasting cluster assignments. In *Proc. of the Conference on Neural Information Processing Systems (NeurIPS)*, 2021.
- [5] N. Chebrolu, P. Lottes, A. Schaefer, W. Winterhalter, W. Burgard, and C. Stachniss. Agricultural Robot Dataset for Plant Classification, Localization and Mapping on Sugar Beet Fields. *Intl. Journal of Robotics Research (IJRR)*, 36:1045–1052, 2017.
- [6] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A Simple Framework for Contrastive Learning of Visual Representations. In *Proc. of the Intl. Conf. on Machine Learning (ICML)*, 2020.
- [7] X. Chen, H. Fan, R. Girshick, and K. He. Improved baselines with momentum contrastive learning. *arXiv preprint*, 2003.04297, 2020.
- [8] E.D. Cubuk, B. Zoph, D. Mane, V. Vasudevan, and Q.V. Le. Autoaugment: Learning augmentation policies from data. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [9] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- [10] D. Erhan, Y. Bengio, A. Courville, P.A. Manzagol, P. Vincent, and S. Bengio. Why does unsupervised pre-training help deep learning? *Journal of Machine Learning Research*, 11(19):625–660, 2010.
- [11] M. Everingham, L. Van Gool, C. Williams, J. Winn, and A. Zisserman. The Pascal Visual Object Classes (VOC) Challenge. *Intl. Journal of Computer Vision (IJCV)*, 88(2):303–338, 2010.
- [12] G. Farjon, Y. Itzhaky, F. Khoroshevsky, and A. Bar-Hillel. Leaf counting: Fusing network components for improved accuracy. *Frontiers in plant science*, 12:575751, 2021.
- [13] F. Fiorani and U. Schurr. Future scenarios for plant phenotyping. *Annual review of plant biology*, 64:267–291, 2013.
- [14] D. Gogoll, P. Lottes, J. Weyler, N. Petrinic, and C. Stachniss. Unsupervised Domain Adaptation for Transferring Plant Classification Systems to New Field Environments, Crops, and Robots. In *Proc. of the IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*, 2020.
- [15] F. Görlich, E. Marks, A.K. Mahlein, K. König, P. Lottes, and C. Stachniss. Uav-based classification of cercospora leaf spot using rgb images. *Drones*, 5(2):34, 2021.
- [16] J.B. Grill, F. Strub, F. Altché, C. Tallec, P.H. Richemond, E. Buchatskaya, C. Doersch, B.A. Pires, Z.D. Guo, M.G. Azar, B. Piot, K. Kavukcuoglu, R. Munos, and M. Valko. Bootstrap your own latent: A new approach to self-supervised learning. In *Proc. of the Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- [17] M. Halstead, A. Ahmadi, C. Smitt, O. Schmittmann, and C. McCool. Crop agnostic monitoring driven by deep learning. *Frontiers in plant science*, 12:786702, 2021.
- [18] Z.K. Hartley and A.P. French. Domain adaptation of synthetic images for wheat head detection. *Plants*, 10(12):2633, 2021.
- [19] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum Contrast for Unsupervised Visual Representation Learning. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [20] K. He, R. Girshick, and P. Dollar. Rethinking ImageNet Pre-training. In *Proc. of the IEEE/CVF Intl. Conf. on Computer Vision (ICCV)*, 2019.
- [21] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [22] M.Y. Liu and O. Tuzel. Coupled generative adversarial networks. In *Proc. of the Conference on Neural Information Processing Systems (NeurIPS)*, 2016.
- [23] J. Long, E. Shelhamer, and T. Darrell. Fully Convolutional Networks for Semantic Segmentation. In *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [24] P. Lottes, J. Behley, N. Chebrolu, A. Milioto, and C. Stachniss. Joint Stem Detection and Crop-Weed Classification for Plant-specific Treatment in Precision Farming. In *Proc. of the IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*, 2018.
- [25] P. Lottes, J. Behley, A. Milioto, and C. Stachniss. Fully convolutional networks with sequential information for robust crop and weed detection in precision farming. *IEEE Robotics and Automation Letters (RA-L)*, 3:3097–3104, 2018.
- [26] P. Lottes and C. Stachniss. Semi-supervised online visual crop and weed classification in precision farming exploiting plant arrangement. In *Proc. of the IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*, 2017.
- [27] F. Magistri, N. Chebrolu, and C. Stachniss. Segmentation-Based 4D Registration of Plants Point Clouds for Phenotyping. In *Proc. of the IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*, 2020.
- [28] E. Marks, F. Magistri, and C. Stachniss. Precise 3D Reconstruction of Plants from UAV Imagery Combining Bundle Adjustment and Template Matching. In *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*, 2022.
- [29] J. McCormac, A. Handa, S. Leutenegger, and A.J. Davison. SceneNet RGB-D: Can 5M Synthetic Images Beat Generic ImageNet Pre-training on Indoor Segmentation? In *Proc. of the IEEE Intl. Conf. on Computer Vision (ICCV)*, 2017.
- [30] A. Milioto, P. Lottes, and C. Stachniss. Real-time Semantic Segmentation of Crop and Weed for Precision Agriculture Robots Leveraging Background Knowledge in CNNs. In *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*, 2018.
- [31] C. Potena, D. Nardi, and A. Pretto. Fast and accurate crop and weed identification with summarized train sets for precision agriculture. In *Proc. of Intl. Conf. on Intelligent Autonomous Systems (IAS)*, 2016.
- [32] A. Pretto, S. Aravecchia, W. Burgard, N. Chebrolu, C. Dornhege, T. Falck, F. Fleckenstein, A. Fontenla, M. Imperoli, R. Khanna, F. Liebisch, P. Lottes, A. Milioto, D. Nardi, S. Nardi, J. Pfeifer, M. Popović, C. Potena, C. Pradalier, E. Rothacker-Feder, I. Sa, A. Schaefer, R. Siegwart, C. Stachniss, A. Walter, W. Winterhalter, X. Wu, and J. Nieto. Building an Aerial-Ground Robotics System for Precision Farming. *IEEE Robotics and Automation Magazine (RAM)*, 28(3):29–49, 2020.
- [33] M.A. Rahman and Y. Wang. Optimizing Intersection-Over-Union in Deep Neural Networks for Image Segmentation. In *Proc. of the Int. Symp. on Visual Computing*, 2016.
- [34] A.J. Ratner, H.R. Ehrenberg, Z. Hussain, J. Dunnmon, and C. Ré. Learning to compose domain-specific transformations for data augmentation. In *Proc. of the Conference on Neural Information Processing Systems (NeurIPS)*, 2017.
- [35] A. Razavian, H. Azizpour, J. Sullivan, and S. Carlsson. CNN Features off-the-shelf: an Astounding Baseline for Recognition. *arXiv preprint*, 1403.6382v3, 2014.
- [36] E. Romera, J.M. Alvarez, L.M. Bergasa, and R. Arroyo. Erfnet: Efficient residual factorized convnet for real-time semantic segmentation. *IEEE Trans. on Intelligent Transportation Systems (ITS)*, 19(1):263–272, 2018.
- [37] E. Tuba, R. Jovanovic, and M. Tuba. Plant diseases detection based on color features and kapur’s method. *World Scientific and Engineering Academy and Society (WSEAS) Trans. Inf. Sci. Appl.*, 14:31–39, 2017.
- [38] O. Vysotska, H. Kuhlmann, and C. Stachniss. UAVs Towards Sustainable Crop Production. In *Workshop at Robotics: Science and Systems*, 2019.
- [39] J. Weyler, A. Milioto, T. Falck, J. Behley, and C. Stachniss. Joint Plant Instance Detection and Leaf Count Estimation for In-Field Plant Phenotyping. *IEEE Robotics and Automation Letters (RA-L)*, 6(2):3599–3606, 2021.
- [40] J. Weyler, J. Quakernack, P. Lottes, J. Behley, and C. Stachniss. Joint Plant and Leaf Instance Segmentation on Field-Scale UAV Imagery. *IEEE Robotics and Automation Letters (RA-L)*, 7(2):3787–3794, 2022.
- [41] J. Weyler, F. Magistri, P. Seitz, J. Behley, and C. Stachniss. In-field phenotyping based on crop leaf and plant instance segmentation. In *Proc. of the IEEE Winter Conf. on Applications of Computer Vision (WACV)*, 2022.
- [42] D.M. Woebbecke, G.E. Meyer, K.V. Bargaen, and D.A. Mortensen. Color indices for weed identification under various soil, residue,

- and lighting conditions. *Transactions of the American Society of Agricultural and Biological Engineers (ASABE)*, 38:259–269, 1994.
- [43] Z. Wu, Y. Xiong, S. Yu, and D. Lin. Unsupervised feature learning via non-parametric instance-level discrimination. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [44] D. Yoo, N. Kim, S. Park, A.S. Paek, and I.S. Kweon. Pixel-level domain transfer. In *Proc. of the Europ. Conf. on Computer Vision (ECCV)*, 2016.
- [45] J. You, W. Liu, and J. Lee. A dnn-based semantic segmentation for detecting weed and crop. *Computers and Electronics in Agriculture*, 178:105750, 2020.
- [46] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny. Barlow twins: Self-supervised learning via redundancy reduction. In *Proc. of the Intl. Conf. on Machine Learning (ICML)*, 2021.
- [47] H. Zhang, M. Cisse, Y.N. Dauphin, and D. Lopez-Paz. mixup: Beyond empirical risk minimization. In *Proc. of the Int. Conf. on Learning Representations (ICLR)*, 2018.
- [48] Z. Zhong, L. Zheng, G. Kang, S. Li, and Y. Yang. Random erasing data augmentation. In *Proc. of the Conference on Advancements of Artificial Intelligence (AAAI)*, 2020.