

# Unsupervised Semantic Label Generation in Agricultural Fields

Gianmarco Roggiolani<sup>1,\*</sup>, Julius Rückin<sup>1</sup>, Marija Popović,<sup>2</sup> Jens Behley,<sup>1</sup> and

Cyrill Stachniss <sup>1,3</sup>

<sup>1</sup>Center for Robotics, University of Bonn, Germany
 <sup>2</sup>Delft University of Technology, Netherlands
 <sup>3</sup>Lamarr Institute for Machine Learning and Artificial Intelligence, Germany

Correspondence\*: Gianmarco Roggiolani groggiol@uni-bonn.de

#### 2 ABSTRACT

1

3 Robust perception systems allow farm robots to recognize weeds and vegetation, enabling the selective application of fertilizers and herbicides to mitigate the environmental impact of traditional 4 5 agricultural practices. Today's perception systems typically rely on deep learning to interpret sensor data for tasks such as distinguishing soil, crops, and weeds. These approaches usually 6 require substantial amounts of manually labeled training data, which is often time-consuming and 7 requires domain expertise. This paper aims to reduce this limitation and propose an automated 8 labeling pipeline for crop-weed semantic image segmentation in managed agricultural fields. It 9 allows the training of deep learning models without or with only limited manual labeling of images. 10 Our system uses RGB images recorded with unmanned aerial or ground robots operating in 11 the field to produce semantic labels exploiting the field row structure for spatially consistent 12 labeling. We use the rows previously detected to identify multiple crop rows, reducing labeling 13 errors and improving consistency. We further reduce labeling errors by assigning an "unknown" 14 class to challenging-to-segment vegetation. We use evidential deep learning because it provides 15 predictions uncertainty estimates that we use to refine and improve our predictions. In this 16 way, the evidential deep learning assigns high uncertainty to the weed class, as it is often 17 less represented in the training data, allowing us to use the uncertainty to correct the semantic 18 predictions. Experimental results suggest that our approach outperforms general-purpose labeling 19 methods applied to crop fields by a large margin and domain-specific approaches on multiple 20 fields and crop species. Using our generated labels to train deep learning models boosts our 21 prediction performance on previously unseen fields with respect to unseen crop species, growth 22 stages, or different lighting conditions. We obtain an IoU of 88.6% on crops, and 22.7% on weeds 23 for a managed field of sugarbeets, where fully supervised methods have 83.4% on crops and 24 33.5% on weeds and other unsupervised domain-specific methods get 54.6% on crops and 25 11.2% on weeds. Finally, our method allows fine-tuning models trained in a fully supervised 26 fashion to improve their performance in unseen field conditions up to +17.6% in mean IoU without 27 additional manual labeling. 28

29 Keywords: Agricultural Automation, Robotic Crop Monitoring, Deep Learning for Agricultural Robots, Semantic Scene Understanding,

30 Automatic Labeling, Unsupervised Learning

## **1 INTRODUCTION**

The demand for food is constantly increasing due to the growing world population, requiring new methods 31 to optimize crop production (Horrigan et al., 2002; Ewert et al., 2023; Storm et al., 2024; Walter et al., 32 33 2017). The use of robotic systems in agriculture has the promise to make processes, such as monitoring fields (Ahmadi et al., 2020; Boatswain Jacques et al., 2021), phenotyping (Weyler et al., 2022b), and 34 35 weed spraying (Wu et al., 2020), more efficient and sustainable (Cheng et al., 2023). Commonly, robotic platforms perceive their environment using deep learning methods to semantically interpret complex data 36 collected with onboard sensors (Dainelli et al., 2024). However, these approaches usually require large 37 amounts of human-labeled data to achieve satisfactory performance for real-world deployment and often 38 fall short in unseen field conditions (Wang et al., 2022; Magistri et al., 2023). 39

In this paper, we examine the problem of automated semantic crop-weed segmentation in RGB images, 40 enabling robots to perform tasks, such as automated weeding (Balabantaray et al., 2024; Saqib et al., 2023), 41 controlled usage of pesticide (Murugan et al., 2020), harvesting (Pan et al., 2023), or phenotyping (Weyler 42 et al., 2022a). We aim to maximize a semantic segmentation neural network's performance in various 43 field deployment conditions, e.g., different growth stages, crop species, or lighting conditions, without 44 human-labeled training data. This is crucial to ensure a robust crop-weed segmentation in new unseen 45 fields to enable robots to perform weeding and harvesting. Our approach automatically labels onboard 46 RGB images based on the robot's pose and the current map of the field semantically segmented into crops 47 and weeds. In this way, semantic labels are generated using the robot's spatial information and the field 48 arrangement's crop row structure. 49

Previous heuristic-based methods for unsupervised semantic segmentation in agriculture proposed by 50 Lottes et al. (2017) and Winterhalter et al. (2018) rely on poorly generalizing assumptions about field 51 52 arrangements, e.g., absence of weeds in the crop row (Lottes et al., 2017), constant distance between plants' rows (Lottes and Stachniss, 2017; Winterhalter et al., 2018), or non-overlapping vegetation 53 54 components (Lottes et al., 2017). Although fully supervised deep learning-based approaches do not rely on geometric assumptions, they rely on in-domain human-labeled training data. The performance of 55 56 such approaches is satisfactory when being deployed in conditions similar to those they were trained on. 57 However, their performance usually rapidly deteriorates in novel deployment conditions, e.g., different crop species, weeds pressure, lighting conditions, or growth stage, requiring new human-labeled training 58 data. This is costly and makes these approaches impractical for application when there is not enough time, 59 60 money, or data to re-train the approach on new field conditions.

The main contribution of this paper is a novel heuristic approach for unsupervised soil-weed-crop 61 segmentation in managed agricultural fields addressing these limitations. Our method automatically 62 generates labels used to train deep semantic segmentation networks. The overview of our pipeline is shown 63 in Fig. 1. Our pipeline takes the current RGB image and the camera pose of the robotic platform as input 64 to compute a semantic map of the field. As a key novelty, we use the semantic map to enforce the spatial 65 consistency of labels. To this end, we propagate the information about the crop rows in the map, leading to 66 better crop segmentation across different growth stages. Additionally, we do not assign labels to vegetation 67 components that are close to the crop rows but are not classified as crops. This reduces possible labeling 68 errors and thus improves model predictions after training on our generated labels. We use the generated 69 image-label pairs to train an uncertainty-aware evidential semantic segmentation network (Sensoy et al., 70 2018). At inference, as a post-processing step, we exploit the predicted uncertainties to refine the final 71 semantic predictions. 72

73 In sum, we make three key claims: our approach (i) generates more accurate semantic labels than previous

unsupervised label generation approaches on multiple crop species, growth stages, and lighting conditions;

75 (ii) we outperform previous unsupervised semantic segmentation approaches by combining our spatially

76 consistent generated labels and uncertainty-aware semantic neural networks; and (iii) improve performance

of fully supervised models on previously unseen crops, growth stages, or soil conditions after fine-tuning

using our automatically generated labels. These claims are backed up by our experimental evaluation. We
 open-source our code upon paper acceptance.

# 79 open-source our code upon paper acceptance.

## 2 RELATED WORK

80 Our work uses heuristic-based computer vision techniques for semantic segmentation of RGB images 81 to automatically generate weed-crop segmentation labels of agricultural fields for training a semantic 82 segmentation network. We train the network in an uncertainty-aware fashion using evidential deep 83 learning (Sensoy et al., 2018) to post-process predictions at inference time based on their uncertainty.

84 Heuristic-based semantic segmentation. Otsu (1979) proposed using gray-level histograms for binary 85 image segmentation based on an automatic threshold assuming a bimodal distribution for fore- and 86 background pixels. Pong et al. (1984) propose the region-growing algorithm segmenting images in multiple 87 regions after providing initial seeds for each region. Similarly, the Watershed algorithm (Najman and 88 Schmitt, 1996) requires user-defined markers to segment objects using a distance function. To overcome 89 the need for initial seeds, Canny (1986) used edge detectors to distinguish regions. To incorporate statistical 90 image features for segmentation, Loyd (1982) adopted the K-means algorithm. To allow automatic 91 robotic intervention in the fields, Riehle et al. (2020) and Gao et al. (2020) applied semantic segmentation 92 techniques to the agricultural domain. Lottes et al. (2017) further advance these general-purpose approaches 93 by exploiting the field arrangement and deploying their method on a ground field robot. Similarly, our approach also exploits the field arrangement to automatically segment images. In contrast, we additionally 94 95 enhance spatial label consistency using robotic semantic mapping. Further, we do not assign labels to image parts likely to include labeling errors. In this way, we reduce the number of erroneous crop and weed 96 instances, which is essential to achieve high prediction performance and consistent uncertainty estimation 97 of the trained deep neural network. 98

99 Learning-based semantic segmentation. Recent approaches mainly use neural networks to extract latent image features for semantic segmentation. Various convolutional neural network 100 101 architectures (Romera et al., 2018; He et al., 2017), and more recently, Vision transformers (Strudel et al., 2021; Cao et al., 2023) have been applied to semantic segmentation. A large portion of these 102 103 approaches have also been evaluated or adapted to the agricultural domain. Cui et al. (2024) propose 104 an improvement to the U-net architecture by Ronneberger et al. (2015) to segment corns and weeds 105 while Zenkl et al. (2022) use the DeepLabV3 architecture by Chen et al. (2017) to segment wheat. These 106 approaches usually require vast amounts of per-pixel human-labeled training data, covering all the desired 107 crop species, growth stages, lighting conditions, and other deployment conditions to ensure promising test-time performance. Hence, many works have investigated how to reduce the labeling effort of deep 108 learning-based approaches. One popular method is pre-training the network on different easy-to-label tasks, 109 e.g., image classification (Deng et al., 2009) or using self-supervision (Chen et al., 2020), and fine-tuning 110 the pre-trained network using few human-labeled per-pixel annotations specific to the target application. 111 Other works propose to train networks on sparse labels instead of dense per-pixel labels (Lee et al., 2022), 112 113 so called weakly supervised semantic segmentation. In the agricultural domain, Zhao et al. (2023) reduce the need for per-pixel labels using scrawl annotations, i.e. manually drawn lines, to weakly supervise a 114

semantic segmentation model. Chen et al. (2024) remove per-pixel annotations completely, only exploiting reference images to localize disease symptoms in plants, using an innovative class activation mapping method. In contrast to Chen et al. (2024), we propose a new unsupervised approach to automatically generate per-pixel semantic segmentation labels exploiting domain knowledge of the field arrangement. Our semantic labels can be directly used for network training without the need for human labels or for fine-tuning pre-trained networks on unseen fields.

Uncertainty-aware deep learning. Classical neural networks are known to often provide 121 overconfident wrong point estimate predictions (Abdar et al., 2021). Several works, including the one 122 by Lakshminarayanan et al. (2017), use ensembles of multiple independently initialized and trained neural 123 networks to quantify predictive uncertainty. Although ensembles improve prediction performance and 124 125 model calibration, they induce high computational costs during training. Gal and Ghahramani (2016) propose Monte Carlo dropout to approximate predictive uncertainty with a single network trained with 126 127 dropout. At inference, multiple forward passes with independently sampled dropout masks are performed to compute predictive uncertainty. Although more compute-efficient at train time, Monte Carlo dropout 128 produces overconfident predictions compared to ensembles (Beluch et al., 2018b). More recently, Sensoy 129 et al. (2018) proposed evidential deep learning for image classification to predict uncertainty using a single 130 forward pass. As evidential deep learning performs on par with ensembles while drastically reducing online 131 compute requirements, we adapt the evidential deep learning framework to our semantic segmentation task 132 using the predictive uncertainties for label post-processing, facilitating deployment on compute-constrained 133 robots. We use the network's uncertainty to correct its prediction about the weeds, which is the most 134 under-represented class and, thus, the most uncertain for the model. 135

Our approach combines a heuristic-based method to automatically generate partial but consistent per-pixel semantic labels. In contrast to learning-based approaches, our approach does not require human-labeled data and, at the same time, improves label consistency and, thus, the network's prediction performance over previous heuristic-based approaches.

## 3 METHOD

We propose a heuristic-based approach to automatically segment RGB images of agricultural 140 fields collected using unmanned ground vehicles (UGVs) or unmanned aerial vehicles (UAVs) in three 141 classes: soil, crop, and weed. Based on the robot's pose, we fuse each generated semantic image label in 142 an online-built global semantic field map. A key aspect of our approach is that we enforce spatial label 143 consistency based on the global semantic field map. To reduce the possibility of labeling errors, we only 144 145 label the detected rows as crops and the vegetation components that are far away from the rows as weeds. In this way, we trade off label quality with quantity to improve prediction performance after training our 146 uncertainty-aware semantic segmentation network (Sensoy et al., 2018) on labels extracted from the global 147 semantic field map. At inference time, we post-process the network's predictions using their associated 148 uncertainty to refine uncertain vegetation predictions. 149

## 150 3.1 Semantic Field Mapping

151 We perform semantic mapping to enforce spatial consistency across automatically generated semantic 152 labels. Furthermore, the semantic map allows us to extract image-label pairs from the map with different 153 rotations, positions, and scales. We assume that our robotic system is equipped with a downwards-facing 154 RGB camera. At each time step t, it collects an image  $I_t \in \mathbb{R}^{H \times W \times 3}$ , where H and W are the height and 155 width of the image, respectively. Let  $\mathbf{p}_t = (x_t, y_t, z_t, \phi_t)^\top$  be the robot pose, where we consider the 3D 156 position  $(x_t, y_t, z_t)$  and the yaw angle  $\phi_t \in (-\pi, \pi]$  with respect to the origin of the mapping mission. 157 Any path is defined by a sequence of poses that we use to fuse our predicted labels in the global semantic 158 field map  $S_t : G \to \mathbb{N}^{K \times \hat{H} \times \hat{W}}$ , where G is a grid discretizing the environment into  $\hat{H} \times \hat{W}$  cells with 159 K possible semantic classes. Each image  $I_t$  along the path is segmented by our approach based on the 160 previous map  $S_{t-1}$  and then fused into the semantic map to compute  $S_t$  accumulating predictions. We use 161 majority voting to assign the most likely class. In practice, we follow a common lawnmower-like coverage 162 path to efficiently cover agricultural fields (Höffmann et al., 2023), as shown in Fig. 2.

#### 163 3.2 Automatic Labeling

At each time step t, our automatic labeling approach takes as input the image  $I_t$  and the semantic field map  $S_{t-1}$  to produce a semantic label for image  $I_t$ . We use the map  $S_{t-1}$  to estimate potential weeds and crops in image  $I_t$  to enforce spatial consistency and reduce labeling errors. Our automatic labeling procedure is exemplarily visualized in Fig. 3 and consists of the following steps: first, we extract the vegetation mask and apply the Hough transform to detect the main crop row in the current image  $I_t$ . Second, we propagate all previously detected lines  $\mathcal{R}_{t-1}$  to the current pose to segment multiple crop rows. Third, we label the vegetation components with a minimal distance to all rows as weeds.

Hough transform. We compute a binary vegetation mask  $I_{t,vm} \in \{0,1\}^{H \times W}$  using graph-based 171 segmentation proposed by Felzenszwalb and Huttenlocher (2004), where a pixel is 1 if it contains vegetation, 172 i.e. crop or weed, and 0 if it contains soil. We apply the Hough transform introduced by Hough (1959) 173 to the vegetation mask  $I_{t,vm}$  to detect crop rows in image  $I_t$ . This gives us a set of supporting lines in  $I_t$ . 174 Each line *i* is parameterized by the distance  $r_{t,i}$  from the image origin to the closest point on the line, and 175 the angle  $\theta_{t,i}$  between the image's x-axis and the line connecting the origin to the closest point on the line. 176 The origin is the lower-left pixel of  $I_t$ . The best-fitting line is the one that maximizes the overlap with the 177 vegetation mask  $I_{t,vm}$ . In Fig. 4, we show an example of a fitted crop row line (white). We discretize the 178 Hough line radius search space using a pixel resolution of  $l_w = 5 \text{ px}$  to robustly fit lines in presence of 179 noisy vegetation masks. We define the minimum number of overlapping pixels  $\tau_{px} = H$  to fit the line along 180 the whole image height. We keep only the best-fitting line of parameters  $(r_t, \theta_t)$  returned from the Hough 181 Transform and add it to the set of the crop rows detected in the map  $\mathcal{R}_t = \mathcal{R}_{t-1} \cup (r_t, \theta_t)$  to use them in 182 the following step. Based on the best-fitting line parameters  $(r_t, \theta_t)$ , we create a binary mask  $I_{t,\text{line}}$ , which 183 is 1 for all pixels on the line and 0 otherwise. We save the line mask to facilitate the computation of the 184 185 following steps. The mask obtained from our example image is shown on top of the vegetation mask in Fig. 3. We transform the line parameters for this time step t into the coordinate system of the mapping 186 mission's origin  $p_0$ . 187

**Propagating predictions.** We use our semantic map  $S_{t-1}$  to retrieve the predicted lines  $\mathcal{R}_{t-1}$  and propagate them into our current image  $\mathbf{I}_t$ . This allows us to predict multiple crop rows consistent with the rows detected in previously explored areas of the crop field. At the first time step t = 0, the semantic map and  $\mathcal{R}_0$  are both empty, thus we skip this step. At each time step  $t \ge 1$ , we compute the position of the newly acquired image in the coordinate system of the initial pose  $\mathbf{p}_0$ , given by the transformation matrix  $\mathbf{T}_t^0 \in \mathbb{R}^{3\times 3}$ . Then, we check which lines in  $\mathcal{R}_{t-1}$  intersect  $\mathbf{I}_t$  and should be propagated into its semantic prediction. For each line i in  $\mathcal{R}_{t-1}$ , we compute the parameters  $r_{t,i}$  and  $\theta_{t,i}$  in the coordinate system of  $\mathbf{p}_0$  as

$$r_{t,i} = \left\| \left( \mathbf{T}_{t}^{0} \right)^{-1} \begin{bmatrix} r_{t-1,i} \cos(\theta_{t-1,i}) \\ r_{t-1,i} \sin(\theta_{t-1,i}) \\ 0 \end{bmatrix} \right\|_{2},$$
(1)

$$\theta_{t,i} = \theta_{t-1,i} - \phi_t,\tag{2}$$

195 where  $r_{t-1,i} \cos(\theta_{t-1,i})$  and  $r_{t-1,i} \sin(\theta_{t-1,i})$  represent the (x, y) coordinates of the closest pixel to the 196 origin for line *i*, assuming flat terrain. We include these lines in  $\mathbf{I}_{t,\text{line}}$ , i.e. we set the pixels covered by 197 these lines to 1. To reduce the computation time, we reject lines that are too close to those already present 198 in the mask  $\mathbf{I}_{t,\text{line}}$ . In particular, we reject line *i* if its distance to any other line in  $\mathcal{R}_t$  is smaller than  $2l_w$ . In 199 Fig. 3, we showcase line propagation from a previous image, enabling us to detect a second crop row on 200 the image's right side.

As we propagate our line predictions from previously recorded images into the current image, we use an eroded version of the vegetation mask  $I_{t,vm}$  to extract single vegetation components. We use a square kernel of size 3 for the erosion to remove noise from  $I_{t,vm}$  and reduce the mislabeling of weeds touching the crops in the crop row. Then, all vegetation components intersecting lines in  $I_{t,line}$  are assigned to the crop class, yielding a new binary mask  $M_t \in \{0, 1\}^{H \times W}$  where a pixel is 1 if it is labeled as crop, and 0 otherwise. We show the result in Fig. 3, where soil is depicted in black and crop is depicted in green. Next, we describe which remaining vegetation components are assigned the weed class.

Weed labeling. Naively classifying any vegetation component in  $I_{t,vm}$  not yet labeled as crop in  $M_t$ usually results in poor weed label quality. Although these remaining vegetation components might be crop, the row detection could have failed because of low sensor resolution, wrong odometry or pose information, or bad lighting conditions (Lottes et al., 2016), such that these crop instances are not included in  $M_t$ . To avoid labeling these potential crops as weeds, we do not label the vegetation components, which are likely to introduce labeling errors and ignore them during network training. To this end, we compute the distance from each of the *N* crop pixels of  $M_t$  with value 1 to their respective closest line as follows

$$d(x,y) = \arg\min_{(r_{t,i},\theta_{t,i})\in\mathcal{R}_t} |x\cos(\theta_{t,i}) + y\sin(\theta_{t,i}) - r_{t,i}|.$$
(3)

We aim to estimate crop sizes along the detected rows using these distances d(x, y). Hence, we use an indicator function  $\mathbb{1}(x, y)$  that returns 1 if the pixel (x, y) is 1 in  $\mathbf{M}_t$  and zero otherwise to extract the mean  $\mu_d = \frac{\sum_{(x,y)} \mathbb{1}(x,y) d(x,y)}{N}$  and standard deviation  $\sigma_d = \sqrt{\frac{\sum_{(x,y)} \mathbb{1}(x,y) (d(x,y) - \mu_d)^2}{N}}$ . We define the minimum distance  $d_{\min}$  required for any unlabeled vegetation instance to be labeled as a weed as

$$d_{\min} = \mu_d + \delta \,\sigma_d,\tag{4}$$

219 where  $\delta = 3$  in our setting, such that only vegetation instances with a large distance from all rows are 220 considered weeds. All vegetation components that were not labeled as crops and whose distance to the 221 lines is smaller than  $d_{\min}$  are left unlabeled. Note that large values of  $\delta$  reduce the number of components 222 labeled as weeds, while small values of  $\delta$  are prone to weed labeling errors. The key idea behind this step is 223 that  $\mu_d$  and  $\sigma_d$  represent the area around the detected crop row where we assume there may be other crops 224 that were not touched by the line and that we leave unlabeled. Outside of this area, we are fairly confident 225 that the vegetation component is a weed as it is far from the detected crop row with plants of estimated size  $\mu_d$ . The resulting label for the example image is shown in Fig. 3, where components close to the crop row on the right are not labeled while the component on the upper-left corner is labeled as a weed.

#### 228 3.3 Learning with Uncertainty

Once we finish our mapping mission as described in Sec. 3.2, we can extract any number of imagelabel pairs with any size, rotation, and aspect ratio. We use the extracted labels to train a semantic segmentation network. We follow the evidential deep learning framework by Sensoy et al. (2018) to predict semantic segmentation and the network's prediction uncertainty at the same time. Estimating the prediction uncertainty allows us to account for the "unknown" class by refining the network's semantic predictions in a post-processing step described in Sec. 3.4.

The key idea behind evidential deep learning is to predict a Dirichlet distribution over all possible class probabilities instead of a single point estimate as in deterministic deep neural networks. In this way, the evidential network minimizes the prediction error while maximizing the prediction uncertainty for ambiguous image parts. We use evidential deep learning instead of Bayesian deep learning approaches (Gal and Ghahramani, 2016; Beluch et al., 2018a) as it is empirically shown to produce similarly or bettercalibrated prediction uncertainties (Sensoy et al., 2018) while being computationally more efficient during training than ensemble methods and during inference than Monte Carlo dropout.

242 We train the network to minimize the Bayes risk cross-entropy for a pixel (x, y) of image I,

$$\mathcal{L}_{CE,(x,y)} = \sum_{k=1}^{K-1} \mathbf{y}_{(x,y),k} \left( \psi(Q_{(x,y)}) - \psi(\boldsymbol{\alpha}_{(x,y),k}) \right),$$
(5)

243 where  $\psi$  is the digamma function,  $\mathbf{y}_{(x,y),k} = 1$  if the pixel (x, y) of I belongs to ground truth class k, 244  $Q_{(x,y)} = \sum_{k=1}^{K} \boldsymbol{\alpha}_{(x,y),k}$ , and  $\boldsymbol{\alpha}_{(x,y),k}$  is the evidence predicted by the network in support of class k. We do 245 not compute this loss for the pixels assigned to the "unknown" class, so we sum only over the remaining 246 K - 1 classes, i.e., soil, crop, and weed. We additionally minimize the Kullback-Leibler (KL) divergence 247 between the uniform  $D(\mathbf{1}_{K-1})$  and predicted Dirichlet distribution  $D(\tilde{\boldsymbol{\alpha}}_{(x,y)})$  for all non-ground-truth 248 classes (Sensoy et al., 2018),

$$\mathcal{L}_{(x,y)} = \mathcal{L}_{CE,(x,y)} + \lambda_{\text{epoch}} KL[D(\tilde{\boldsymbol{\alpha}}_{(x,y)})||D(\mathbf{1}_{K-1})],$$
(6)

249

$$\tilde{\boldsymbol{\alpha}}_{(x,y),k} = \mathbf{y}_{(x,y),k} + (1 - \mathbf{y}_{(x,y),k})\boldsymbol{\alpha}_{(x,y),k},\tag{7}$$

for all K - 1 classes, and  $\lambda_{epoch} = \min(1.0, \frac{epoch}{T})$  with epoch being the current training epoch and T the maximum annealing epoch. We minimize the overall training loss

$$\mathcal{L} = \frac{1}{HW} \sum_{x=1}^{H} \sum_{y=1}^{W} \mathcal{L}_{(x,y)},\tag{8}$$

which is the average over all image pixels, iterating over all training images. At inference time, the network predicts the semantic class and an uncertainty for each pixel, that we use for our label refinement.

#### 254 3.4 Uncertainty-based Label Refinement

We use the network's predicted Dirichlet distribution  $D(\alpha_{(x,y)})$  over all K - 1 classes to quantify the prediction uncertainty for post-processing and refining the predicted semantic labels. The network's prediction uncertainty (Sensoy et al., 2018) for a pixel (x, y) of image I is given by

$$u_{t,(x,y)} = \frac{K-1}{\sum_{k=1}^{K-1} \alpha_{(x,y),k}},$$
(9)

where K-1 is the number of classes without the "unknown" class. In our crop-weed segmentation case, the most under-represented class is weed. Thus, the network will be more uncertain about areas representing weeds than the other classes. We define an adaptive threshold to select the most uncertain pixels (x, y) in any image I as

$$\tau = \frac{\max(u_{(x,y)}) - \min(u_{(x,y)})}{2} + \min(u_{(x,y)}).$$
(10)

We compute a binary mask  $\mathbf{U}_t \in \{0,1\}^{H \times W}$  where a pixel (x, y) is 1, if  $u_{(x,y)} > \tau$ , and 0 otherwise. 262 We compute the connected components of our semantic prediction, aiming to use the ratio between 263 the size of the object and its number of uncertain pixels to refine the component's label. Most of the 264 vegetation components have high uncertainty at their instance boundaries. Instead, we are interested 265 in those components for which also large amounts of interior pixels are uncertain. We iterate over all 266  $c \in \{1, \ldots, C\}$  crop components in our network's prediction and compute for each one a binary mask 267  $\mathbf{C}_c \in \{0, 1\}^{H \times W}$ , which is 1 for all pixels belonging to the component. We also compute their bounding box  $\mathbf{b}_c = (b_c^x, b_c^y, b_c^{\text{height}}, b_c^{\text{width}})$ , where  $b_c^x$  and  $b_c^y$  are the coordinates of the upper left corner of the bounding box, while  $b_c^{\text{height}}$  and  $b_c^{\text{width}}$  are the height and width of the bounding box. We define an adaptive threshold 268 269 270

$$\tau_c = \frac{1}{4} \min\left(\frac{b_c^{\text{width}}}{b_c^{\text{height}}}, \frac{b_c^{\text{height}}}{b_c^{\text{width}}}\right).$$
(11)

This threshold helps us avoid detecting as weeds a lot of small spikes of uncertainty that could arise because of shadows, reflections, or insects. In this way, we only act upon vegetation components where there is a large uncertain area. If the network is uncertain about the prediction of crop component c, it holds that

$$\frac{\sum_{(x,y)} \mathbf{U}_{(x,y)} \mathbf{C}_{c,(x,y)}}{b_c^{\text{width}} b_c^{\text{height}}} > \tau_c.$$
(12)

275 If crop component *c* fullfills Eq. 12, we assign the component's uncertain pixels (x, y) with  $U_{(x,y)} = 1$ 276 to the weed class. We do not re-assign the whole vegetation component as a weed because our network 277 does not provide instances. Hence, there may be components that contain both weeds and crops. These 278 components likely have higher uncertainty since they are labeled as "unknown" and thus being ignored 279 during training. We show in Fig. 5 the result of our post-processing operation for an example image, 280 highlighting the correspondence between the network's wrong predictions, the estimated uncertainty and 281 the post-processed semantic prediction.

# 4 **RESULTS**

The main focus of this work is an automatic labeling pipeline for semantic soil-weed-crop segmentation of RGB images. The results of our experiments support our key claims: our approach (i) generates more accurate semantic labels than previous unsupervised label generation approaches on multiple datasets; (ii) we outperform previous unsupervised semantic segmentation approaches by combing our spatially consistent generated labels and uncertainty-aware semantic neural networks; and (iii) we improve the performance of fully supervised models on previously unseen crops, growth stages, and soil conditions after fine-tuning the network using our automatically generated labels.

#### 289 4.1 Experimental Setup

290 **Datasets.** We use four datasets, three of which are publicly available: PhenoBench (Weyler et al., 2024), as well as the Carrots and Onions from Lincoln University (Bosilj et al., 2020), and a Sugar Beets 291 dataset introduced by Weyler et al. (2022b). The Carrots dataset was recorded in Lincolnshire, UK, in June. 292 The field is under substantial weed pressure and contains weeds with a similar appearance to the crop. 293 294 Furthermore, several regions of vegetation contain crops and weeds in close proximity. The Onions dataset was also recorded in Lincolnshire, UK, but in April. The weed pressure is lower compared to the Carrots 295 296 dataset. The PhenoBench dataset was recorded in Meckenheim, Germany, on different dates between May 297 and June to capture different growth stages. The field contains two varieties of sugar beets and six different 298 weed varieties. The weed pressure varies as the dataset contains images from fields that were fully, partially, or not treated at all with herbicides. The Sugar Beets dataset was also recorded in Meckenheim, Germany, 299 over five different weekly sessions. The field is arranged with small spacing between plants and shows 300 high weed pressure, inducing challenging conditions. We refer to Tab. 1 for information about the camera, 301 image resolution, and ground sampling distance of the datasets. 302

303 Training Details and Hyperparameters. We use ERFNet (Romera et al., 2018) as our network trained 304 using the Adam (Kingma and Ba, 2015) optimizer, a learning rate of 0.01, and a batch size of 32. We set T = 25 in Eq. 6 to linearly increase  $\lambda_{epoch}$  over the first 25 epochs. We report all the hyperparameters 305 306 of our method with their values in Tab. 2. To evaluate the quality of the labels, we generate labels for 307 the validation sets of PhenoBench and Sugar Beets, as well as for the whole Carrots and Onions dataset. Second, we automatically generate labels for the images in their training sets to train our network and 308 evaluate the results on the manually annotated validation sets. We do not split Carrots and Onions to train 309 on them since they consist of only 20 images each. Thus, we do not use them for model training. Instead, 310 we evaluate our label generation and the generalization capabilities of fine-tuned models on these datasets. 311

Metrics. We use the intersection over union (IoU) (Everingham et al., 2010) as a metric for all of our experiments. For the automatic labeling pipeline, we also report the boundary IoU (Cheng et al., 2021) to have a better understanding of the approaches' limitations. The reported mean IoU (mIoU) values are the macro-averages over all classes.

**Baselines.** We use three baselines: two are general-purpose unsupervised semantic segmentation networks not specifically developed for the agricultural domain, while one is an automatic labeling method specifically developed for the agricultural domain. The first baseline is STEGO by Hamilton et al. (2022), which provides an official implementation for the evaluation alongside their models. We use the model trained on MS COCO (Lin et al., 2014) with the vision transformer architecture (Dosovitskiy et al., 2021). STEGO predicts different per-pixel features and then clusters them using self and cross attention mechanisms (Vaswani et al., 2017). Our second baseline is U2Seg by Niu et al. (2024), which builds on

top of STEGO and uses instance information to overcome some of the limitations of the previous work; 323 324 they also open-source their code and provide their models. U2Seg proposes a universal segmentation, coupling instances and semantic classes at training time, to predict clusters at inference time for which 325 they recover class and instance labels. We use the model trained on Imagenet (Deng et al., 2009) and MS 326 COCO with 800 clusters. Lottes et al. (2016) propose a domain-specific method for generating per-pixel 327 crop and weed labels. They use a vegetation mask to detect the main crop row and then label all other 328 vegetation components as weeds. We use their official implementation, removing the NIR image channels. 329 We evaluate their automatically generated labels (base) and the performance of ERFNet trained on their 330 labels (learned). We train the same network with the same training hyperparameters on their and our 331 generated label to ensure a fair comparison. We report the results of ERFNet trained on the manually 332 annotated training set of PhenoBench and evaluated on the validation set as an upper performance bound. 333

#### 334 4.2 Automatic Labeling

In the first experiment, we show that our automatic labeling pipeline generates more accurate semantic soil-weed-crop labels than other methods on multiple datasets. We compare against two general-purpose unsupervised semantic segmentation networks and the domain-specific approach by Lottes et al. (2016).

We show the results on all four datasets in Tab. 3. The general-purpose approaches perform worse than 338 339 the domain-specific methods across all datasets, except for U2Seg on the Onions dataset. As Onions have thin leaves, they are hard to detect with common color histogram thresholding methods, such as the one by 340 Lottes et al. (2016). Furthermore, the weeds in this dataset are the same size as the crops, leading to bias in 341 crop row detection and introducing a higher risk of confusing weeds and crops. Our approach for label 342 generation, referred to as Ours (base), shows higher crop label quality than Lottes (base) while performing 343 on par or better in terms of weed label quality. Particularly, Lottes (base) confuses substantially more weeds 344 with crops, while our approach, by design, does not assign labels to hard-to-label vegetation components, 345 as described in Sec. 3.2. The Carrots dataset is the only one where U2Seg outperforms the domain-specific 346 approaches, which suffer from the weed pressure when estimating the crop rows. Our method consistently 347 outperforms all other baselines across all datasets with different crop species, weed pressure, growth stages 348 and lighting conditions. Most approaches fail on the Onions dataset due to brighter illumination and thin 349 crops. In contrast, our approach improves by approx. 9% mIoU over the second-best baseline, U2Seg, 350 importantly showing highest improvements in both vegetation classes. 351

352 The boundary IoU confirms the result of the standard IoU metric. As shown in Tab. 3, the approach by Lottes et al. (2016) poorly segments boundaries on most of the datasets. This might be due to wrongly 353 segmented vegetation masks. Aiming to include the boundary of weeds more accurately may worsen the 354 overall performance since soil could be wrongly considered as vegetation. We hypothesise that our approach 355 might suffer from the same problem on the Carrots dataset. The difference between IoU and boundary 356 IoU per class suggests that we underestimate the size of weeds, i.e., high IoU but low boundary IoU for 357 weeds, and overestimate crop size, i.e., low IoU but high boundary IoU for crops. On the Carrots dataset 358 U2Seg outperforms the other methods on the weeds boundary IoU. The weed IoU suggests that U2Seg 359 overestimates weeds, thus obtaining a boost as the total number of pixels in the IoU computation is low. On 360 the Onions dataset, our method's IoU and boundary IoU are almost the same irrespective of the semantic 361 class since the crops and weeds are thin. Thus, the boundary area covers the whole vegetation instance. The 362 other approaches fail to correctly assign weed and crop boundaries on the Onion dataset, which follows 363 from the low weed and crop IoU. On the Sugar Beets dataset, all approaches fail to predict boundaries, most 364 likely due to unusually high weed pressure. Our method accurately segments soil boundaries, suggesting 365

that it at least successfully differentiates between soil and vegetation. Overall, the results suggest that most approaches underestimate the size of vegetation, both crops and weeds. Instead, our automatic labeling method shows the strongest boundary segmentation performance across all methods and classes on most datasets, often by a large margin compared to the second-best method. This further verifies our claim that our automatic labeling pipeline generates more accurate semantic soil-weed-crop labels than previous methods. We show qualitative results of Lottes et al. (2016) and our approach in Fig. 6.

#### 372 4.3 Unsupervised Semantic Segmentation

373 The second experiment evaluates the performance of our automatic label generation combined with 374 network training and uncertainty post-processing on the PhenoBench dataset. We show that training the evidential ERFNet using our automatically generated labels outperforms other unsupervised semantic 375 segmentation models. The general-purpose learning-based approaches have not been fine-tuned on human-376 377 labeled field images to ensure a fair comparison. Our approach and Lottes et al. (2016) generate labels on the PhenoBench training set. We use the public training set of images to have a fair comparison with the 378 fully suprvised ERFNet model, trained on the manual labels. Trained models are evaluated on the official 379 PhenoBench validation set. 380

381 Tab. 4 summarizes the results. We use (learned) to refer to the results obtained by ERFNet after being trained on the labels generated by the approach, and we use (+ uncertainty) to refer to the previous results 382 once we post-process them using the uncertainty estimated by the model. The approach by Lottes et al. 383 384 (2016) confuses more crops with weeds since it naively assigns all vegetation components that are not on the main crop row to the weed class. Hence, Lottes et al. (2016) introduce inconsistent labels in the model's 385 training data. Thus, training the ERFNet on Lottes' labels does not yield uncertainty estimations that are 386 useful for improving the predictions during post-processing. Additionally, this leads to smaller performance 387 improvements after training on their labels than after training on our labels. Using our generated labels to 388 train the ERFNet substantially improves the weed and crop predictions over directly using our generated 389 labels. We further improve mIoU and weed predictions by exploiting the estimated uncertainties in Ours 390 391 (uncertainty) for post-processing. Most importantly, Ours (uncertainty) noticeably closes the performance gap between fully supervised and state-of-the-art unsupervised approaches. However, the ERFNet trained 392 on human-labeled training images still predicts weeds more accurately. As the fully supervised model 393 predicts more weeds, it also confuses weeds with crops more often. Hence, our approach performs better 394 395 on both the crop and soil classes. This experiment confirms that our method's conservative approach to labeling, ignoring vegetation components likely to introduce labeling errors combined with evidential deep 396 learning, is a viable solution to largely reduce the need for manually annotated images. 397

#### 398 4.4 Generalization Capability

In the third experiment, we show that our approach enhances the performance of networks trained in a 399 fully supervised fashion by fine-tuning on unseen fields using our automatically generated labels. We do 400 401 not use our evidential network but train an ERFNet using the standard cross-entropy loss to seamlessly 402 fine-tune existing networks pre-trained in a fully supervised fashion. We train two ERFNets, one on each of the human-labeled training sets of PhenoBench and Sugar Beets. We deploy the two models on all four 403 404 datasets without fine-tuning. Then, we fine-tune the two models leveraging our automatically generated 405 labels for the Sugar Beets and PhenoBench datasets. Each model is fine-tuned on the dataset that it was not trained on. 406

407 In Tab. 5, we show the performance of the two models. In brackets, we provide the performance difference 408 after fine-tuning, where blue and red indicate performance improvements or degradations, respectively. The gray rows show the models' performances on the dataset they were trained on. Due to the domain 409 gap between datasets, the models that were not fine-tuned have a lower performance when evaluated on 410 unseen data. Fine-tuning the models makes the performance over the original training data worse as they 411 aim to learn features that are common to both datasets. Our results suggest that using our automatically 412 413 generated labels helps to close the performance gap on previously unseen datasets with different crops, soil 414 types, lighting conditions, and sensor setups. Generally, our fine-tuned models perform better on all classes and datasets, even on the Onions and Carrots datasets, the model was not pre-trained nor fine-tuned on. 415 Only the model that is fine-tuned on the Sugar Beets dataset does not improve performance on the Carrots 416 dataset. We hypothesize this is because the PhenoBench dataset is approx.  $10 \times$  larger than Sugar Beets 417 introducing data imbalance while automatically generated Sugar Beets labels are of lower quality than 418 labels generated on PhenoBench. In sum, using our automatically generated labels helps to fine-tune fully 419 supervised models, enabling better adaptation to unseen field conditions without any additional human 420 labeling costs. 421

# 5 DISCUSSION

A robust perception system is crucial for the successful deployment of robotic platforms in arable fields. 422 Most perception systems rely on data-driven machine learning approaches to train vision models that 423 424 automatically interpret the data collected with onboard sensors, such as RGB cameras. Thus, reliable 425 and accurate learning-based perception systems are crucial to providing valuable information to farmers or autonomous robots. Most learning-based semantic segmentation approaches assume access to large 426 amounts of human-labeled data required to train the vision model. However, their performance rapidly 427 decreases in field conditions they were not trained on, i.e., different crop species, growth stages, weed 428 pressure, and lighting conditions. 429

430 To address this issue, we proposed an automatic labeling approach to obtain semantic information from RGB images of agricultural fields. Our method shows semantic segmentation performance close to the 431 performance of a model trained on large amounts of human-labeled data in a fully supervised fashion. This 432 significantly reduces the need for manually annotated data, reducing costs and relaxing the need for domain 433 experts. The arable field dataset works considered in our experimental evaluation report an average of 2 434 hours per image for labeling the Onions dataset, 3-4 hours per image for the Carrots dataset, and 1-3.5 435 hours *per image* for the PhenoBench dataset. All of the datasets went through at least two labeling rounds, 436 doubling the costs. This highlights the need for new labeling paradigms beyond fully supervised model 437 438 training while maintaining strong prediction performance. Our method is a crucial step towards closing the 439 performance gap between models trained in an unsupervised fashion and fully supervised models without adding additional labeling costs. 440

In our experiments, we show that the fully supervised approach has a lower performance in segmenting 441 crops compared to our unsupervised method, as it is trained on more weed instances. Nevertheless, the 442 fully supervised method still shows the highest mIoU. The unsupervised methods are not exposed to 443 enough weed labels, making them assign the crop class more often. Since the number of crop pixels is 444 generally higher, these errors have a smaller impact on the crop than on the weed segmentation. We also 445 investigate how to use our automatic labeling in combination with supervised methods to improve the 446 447 overall performance in challenging scenarios, i.e., in unseen fields with new crop species and different weed pressure. Fine-tuning comes at the cost of performing worse on the pre-training dataset, as shown in 448

Tab. 5. The degradation largely depends on the size and similarity of the pre-training and automatically labeled dataset used for fine-tuning. Future work could investigate continuous learning methods to train on the newly automatically labeled images without catastrophically forgetting what has already been learned.

452 The need for posed images can be a limitation of our method as it cannot be applied to a dataset of 453 unposed images. However, most of the agricultural datasets are recorded using aerial or ground vehicles 454 that, by default, provide spatial information while recording images in the field, often using GNSS systems 455 such as GPS. Furthermore, we assume deployment in a managed agricultural field. If this assumption 456 does not hold and the weeds are larger than the crops, our crop row detection fails and leads to degraded 457 results. Our results, as well as those by Lottes et al. (2016), show that we could make use of a better 458 vegetation mask to improve unsupervised methods. One possible solution would be to use NIR images, 459 which are less dependent on the lighting conditions compared to RGB images. NIR images are already 460 commonly used for crop segmentation in agriculture (Sahin et al., 2023; Colorado et al., 2020). Moreover, our approach leverages uncertainty estimates to post-process semantic predictions. Current state-of-the-art 461 methods are known to produce partially miscalibrated uncertainty estimates (Beluch et al., 2018a). Thus, 462 our post-processing could benefit from improvements in uncertainty-aware deep learning. Finally, we plan 463 to deploy our approach on a real robot to perform field trials. 464

# 6 CONCLUSION

465 In this paper, we presented a novel approach to automatically generate semantic soil-crop-weed labels of 466 images from agricultural fields. We evaluated our approach on four datasets recorded with different robotic 467 platforms and in various fields. Our approach outperforms previous domain-agnostic and domain-specific 468 unsupervised labeling approaches. Furthermore, we showed that our generated labels allow fine-tuning 469 networks trained in a fully supervised fashion on one dataset to other agricultural fields, e.g., different species, growth stages, and field conditions. In this way, our approach increases the semantic segmentation 470 471 generalization capabilities of existing networks for soil-weed-crop segmentation without additional human 472 labeling effort.

## CONFLICT OF INTEREST STATEMENT

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

## **AUTHOR CONTRIBUTIONS**

475 The authors confirm contribution to the paper as follows: study conception and design: G. Roggiolani, J.

476 Rückin; analysis and interpretation of results: G. Roggiolani; draft manuscript preparation: G. Roggiolani,

477 J. Rückin, M. Popovic, J. Behley, and C. Stachniss; funding acquisition and project coordination: C.

478 Stachniss. All authors reviewed the results and approved the final version of the manuscript.

## FUNDING

This work has been funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation)
under Germany's Excellence Strategy, EXC2070-390732324-PhenoRob.

#### DATA AVAILABILITY STATEMENT

- 481 The publicly available datasets can be found at (PhenoBench) https://www.phenobench.org/dataset.html,
- 482 (Onions) https://lcas.lincoln.ac.uk/nextcloud/index.php/s/RYni5xngnEZEFkR, and (Carrots)
- 483 https://lcas.lincoln.ac.uk/nextcloud/index.php/s/e8uiyrogObAPtcN.

## REFERENCES

- Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., et al. (2021). A review
   of uncertainty quantification in deep learning: Techniques, applications and challenges. *Information fusion* 76, 243–297
- Ahmadi, A., Nardi, L., Chebrolu, N., and Stachniss, C. (2020). Visual Servoing-based Navigation for
  Monitoring Row-Crop Fields. In *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*
- Balabantaray, A., Behera, S., Liew, C., Chamara, N., Singh, M., Jhala, A. J., et al. (2024). Targeted weed
  management of palmer amaranth using robotics and deep learning (yolov7). *Frontiers in Robotics and AI* 11. doi:10.3389/frobt.2024.1441371
- 492 Beluch, W. H., Genewein, T., Nurnberger, A., and Kohler, J. M. (2018a). The power of ensembles for
- active learning in image classification. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*. 9368–9377
- Beluch, W. H., Genewein, T., Nürnberger, A., and Köhler, J. M. (2018b). The power of ensembles for
  active learning in image classification. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*
- Boatswain Jacques, A. A., Adamchuk, V. I., Park, J., Cloutier, G., Clark, J. J., and Miller, C. (2021).
  Towards a machine vision-based yield monitor for the counting and quality mapping of shallots. *Frontiers in Robotics and AI* 8. doi:10.3389/frobt.2021.627067
- Bosilj, P., Aptoula, E., Duckett, T., and Cielniak, G. (2020). Transfer learning between crop types for
  semantic segmentation of crops versus weeds in precision agriculture. *Journal of Field Robotics (JFR)*37, 7–19
- Canny, J. F. (1986). A computational approach to edge detection. *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)* 8, 679–698. doi:10.1109/TPAMI.1986.4767851
- Cao, H., Wang, Y., Chen, J., Jiang, D., Zhang, X., Tian, Q., et al. (2023). Swin-unet: Unet-like pure
   transformer for medical image segmentation. In *Proc. of the Europ. Conf. on Computer Vision (ECCV)*
- Chen, J., Guo, J., Zhang, H., Liang, Z., and Wang, S. (2024). Weakly supervised localization model for
   plant disease based on siamese networks. *Frontiers in Plant Science* 15. doi:10.3389/fpls.2024.1418201
- 510 Chen, L., Papandreou, G., Schroff, F., and Adam, H. (2017). Rethinking Atrous Convolution for Semantic
  511 Image Segmentation. *arXiv preprint* arXiv:1706.05587
- 512 Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. (2020). A Simple Framework for Contrastive Learning
  513 of Visual Representations. In *Proc. of the Intl. Conf. on Machine Learning (ICML)*
- Cheng, B., Girshick, R., Dollár, P., Berg, A. C., and Kirillov, A. (2021). Boundary iou: Improving
  object-centric image segmentation evaluation. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*. 15334–15342
- 517 Cheng, C., Fu, J., Su, H., and Ren, L. (2023). Recent advancements in agriculture robots: Benefits and
   518 challenges. *Machines* 11, 48. doi:10.3390/machines11010048
- 519 Colorado, J. D., Calderon, F. C., Mendez, D., Petro, E., Rojas, J. P., Correa, E. S., et al. (2020). A novel
- 520 nir-image segmentation method for the precise estimation of above-ground biomass in rice crops. *PLOS*
- *ONE* 15, e0239591. doi:10.1371/journal.pone.0239591

- Cui, J., Tan, F., Bai, N., and Fu, Y. (2024). Improving u-net network for semantic segmentation of corns and
   weeds during corn seedling stage in field. *Frontiers in Plant Science* 15. doi:10.3389/fpls.2024.1344958
- Dainelli, R., Bruno, A., Martinelli, M., Moroni, D., Rocchi, L., Morelli, S., et al. (2024). Granoscan: an
  ai-powered mobile app for in-field identification of biotic threats of wheat. *Frontiers in Plant Science* 15.
  doi:10.3389/fpls.2024.1298791
- 527 Deng, J., Dong, W., Socher, R., Li, L., Li, K., and Fei-Fei, L. (2009). Imagenet: A large-scale hierarchical
  528 image database. In *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., et al. (2021). An
   Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *Proc. of the Intl. Conf. on Learning Representations (ICLR)*
- Everingham, M., Van Gool, L., Williams, C., Winn, J., and Zisserman, A. (2010). The Pascal Visual Object
  Classes (VOC) Challenge. *Intl. Journal of Computer Vision (IJCV)* 88, 303–338
- Ewert, F., Baatz, R., and Finger, R. (2023). Agroecology for a sustainable agriculture and food system:
  From local solutions to large-scale adoption. *Annual Review of Resource Economics* 15, 351–381.
  doi:10.1146/annurev-resource-102422-090105
- 537 Felzenszwalb, P. F. and Huttenlocher, D. P. (2004). Efficient graph-based image segmentation.
  538 *Intl. Journal of Computer Vision (IJCV)* 59, 167–181
- Gal, Y. and Ghahramani, Z. (2016). Dropout as a bayesian approximation: Representing model uncertainty
  in deep learning. In *Proc. of the Intl. Conf. on Machine Learning (ICML)*
- Gao, J., Wang, B., Wang, Z., Wang, Y., and Kong, F. (2020). A wavelet transform-based image segmentation
  method. *Intl. Journal for Light and Electron Optics* 208, 164123. doi:https://doi.org/10.1016/j.ijleo.
  2019.164123
- Hamilton, M., Zhang, Z., Hariharan, B., Snavely, N., and Freeman, W. T. (2022). Unsupervised
  semantic segmentation by distilling feature correspondences. In *Proc. of the Intl. Conf. on Learning Representations (ICLR)*
- He, K., Gkioxari, G., Dollár, P., and Girshick, R. (2017). Mask R-CNN. In *Proc. of the IEEE Intl. Conf. on Computer Vision (ICCV)*
- Horrigan, L., Lawrence, R. S., and Walker, P. (2002). How sustainable agriculture can address the
  environmental and human health harms of industrial agriculture. *Environmental health perspectives* 110,
  445–456. doi:10.1289/ehp.02110445
- Hough, P. V. C. (1959). Machine analysis of bubble chamber pictures. In *Proc. of the Intl. Conf. on High-Energy Accelerators and Instrumentation*
- Höffmann, M., Patel, S., and Büskens, C. (2023). Optimal coverage path planning for agricultural vehicles
  with curvature constraints. *Agriculture* 13, 2112. doi:10.3390/agriculture13112112
- Kingma, D. and Ba, J. (2015). Adam: A Method for Stochastic Optimization. In *Proc. of the Intl. Conf. on Learning Representations (ICLR)*
- Lakshminarayanan, B., Pritzel, A., and Blundell, C. (2017). Simple and scalable predictive uncertainty
   estimation using deep ensembles. In *Proc. of the Conf. Neural Information Processing Systems (NIPS)*
- Lee, J., Oh, S. J., Yun, S., Choe, J., Kim, E., and Yoon, S. (2022). Weakly supervised semantic
   segmentation using out-of-distribution data. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*
- Lin, T., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., et al. (2014). Microsoft COCO:
  Common Objects in Context. In *Proc. of the Europ. Conf. on Computer Vision (ECCV)*

565 566 567	Lottes, P., Höferlin, M., Sander, S., Müter, M., Schulze-Lammers, P., and Stachniss, C. (2016). An Effective Classification System for Separating Sugar Beets and Weeds for Precision Farming Applications. In <i>Proc. of the IEEE Intl. Conf. on Robotics &amp; Automation (ICRA)</i>
568	Lottes, P., Höferlin, M., Sander, S., and Stachniss, C. (2017). Effective Vision-based Classification
569	for Separating Sugar Beets and Weeds for Precision Farming. Journal of Field Robotics (JFR) 34,
570	1160–1178
571	Lottes P and Stachniss C (2017) Semi-supervised online visual crop and weed classification in precision
572	farming exploiting plant arrangement. In Proc. of the IEEE/RSI Intl. Conf. on Intelligent Robots and
573	Systems (IROS)
575	Loud S. D. (1092). Loost assume substitution in nome IEEE Trans. on Liferon stien Theory 29, 120, 127.
574 575	doi:10.1109/TIT.1982.1056489
576	Magistri, F., Weyler, J., Gogoll, D., Lottes, P., Behley, J., Petrinic, N., et al. (2023). From one field to
577	another – unsupervised domain adaptation for semantic segmentation in agricultural robotics. Computers
578	and Electronics in Agriculture 212, 108114. doi:https://doi.org/10.1016/j.compag.2023.108114
579	Murugan, K., Shankar, B. J., Sumanth, A., Sudharshan, C. V., and Reddy, G. V. (2020). Smart automated
580	pesticide spraying bot. In Proc. of the Intl. Conf. on Intelligent Sustainable Systems (ICISS)
581	Najman, L. and Schmitt, M. (1996). Geodesic saliency of watershed contours and hierarchical segmentation.
582	IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI) 18, 1163–1173. doi:10.1109/34.
583	546254
584	Niu D. Wang X. Han, X. Lian, L. Herzig R. and Darrell T. (2024) Unsupervised universal image
585	segmentation. In Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)
586	Otsu N (1979) A threshold selection method from gray-level histograms <i>IFFF Trans on Systems Man</i>
587	and Cybernetics 9, 62–66
588	Pan Y Magistri E Läbe T Marks E Smitt C McCool C et al (2023) Panontic mapping with
589	fruit completion and pose estimation for horticultural robots. In <i>Proc. of the IFFE/RSI Intl. Conf. on</i>
590	Intelligent Robots and Systems (IROS)
591	Pong T-C Shapiro I G Watson I T and Haralick R M (1984) Experiments in segmentation using
592	a facet model region grower. <i>Computer Vision, Graphics, and Image Processing</i> 25, 1–23
593	Riehle, D., Reiser, D., and Griepentrog, H. W. (2020). Robust index-based semantic plant/background
594	segmentation for rgb- images. Computers and Electronics in Agriculture 169, 105–201. doi:https:
595	//doi.org/10.1016/j.compag.2019.105201
596	Romera, E., Alvarez, J. M., Bergasa, L. M., and Arroyo, R. (2018). ERFNet: Efficient Residual Factorized
597	ConvNet for Real-Time Semantic Segmentation. IEEE Trans. on Intelligent Transportation Systems
598	( <i>TITS</i> ) 19, 263–272
599	Ronneberger, O., P.Fischer, and Brox, T. (2015). U-Net: Convolutional Networks for Biomedical Image
600	Segmentation. In Proc. of the Medical Image Computing and Computer-Assisted Intervention (MICCAI)
601	Sahin, H. M., Miftahushudur, T., Grieve, B., and Yin, H. (2023). Segmentation of weeds and crops using
602	multispectral imaging and crf-enhanced u-net. Computers and Electronics in Agriculture 211, 107956.
603	doi:https://doi.org/10.1016/j.compag.2023.107956
604	Saqib, M. A., Aqib, M., Tahir, M. N., and Hafeez, Y. (2023). Towards deep learning based smart farming for
605	intelligent weeds management in crops. Frontiers in Plant Science 14. doi:10.3389/fpls.2023.1211235
606	Sensoy, M., Kaplan, L., and Kandemir, M. (2018). Evidential deep learning to quantify classification
607	uncertainty. In Proc. of the Conf. on Neural Information Processing Systems (NeurIPS)

- Storm, H., Seidel, S., Klingbeil, L., Ewert, F., Vereecken, H., Amelung, W., et al. (2024). Research
  Priorities to Leverage Smart Digital Technologies for Sustainable Crop Production. *European Journal of Agronomy* 156, 127178. doi:https://doi.org/10.1016/j.eja.2024.127178
- 611 Strudel, R., Garcia, R., Laptev, I., and Schmid, C. (2021). Segmenter: Transformer for semantic 612 segmentation. In *Proc. of the IEEE/CVF Intl. Conf. on Computer Vision (ICCV)*
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., et al. (2017). Attention Is All
  You Need. In *Proc. of the Conf. on Neural Information Processing Systems (NeurIPS)*
- Walter, A., Finger, R., Huber, R., and Buchmann, N. (2017). Smart farming is key to developing sustainable
  agriculture. *Proceedings of the National Academy of Sciences* 114, 6148–6150. doi:1/10.1073/pnas.
  1707462114
- Wang, J., Lan, C., Liu, C., Ouyang, Y., Qin, T., Lu, W., et al. (2022). Generalizing to unseen domains: A
  survey on domain generalization. *IEEE Trans. on Knowledge and Data Engineering* 35, 8052–8072.
  doi:10.1109/TKDE.2022.3178128
- Weyler, J., Magistri, F., Seitz, P., Behley, J., and Stachniss, C. (2022a). In-Field Phenotyping Based
  on Crop Leaf and Plant Instance Segmentation. In *Proc. of the IEEE Winter Conf. on Applications of Computer Vision (WACV)*
- *Computer Vision (WACV)*Weyler, J., Magistri, F., Marks, E., Chong, Y. L., Sodano, M., Roggiolani, G., et al. (2024). Phenobench:
  A large dataset and benchmarks for semantic image interpretation in the agricultural domain. *IEEE*
- *Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, 1–12doi:10.1109/TPAMI.2024.3419548
  Weyler, J., Quakernack, J., Lottes, P., Behley, J., and Stachniss, C. (2022b). Joint plant and leaf instance
- segmentation on field-scale uav imagery. *IEEE Robotics and Automation Letters (RA-L)* 7, 3787–3794
  Winterhalter, W., Fleckenstein, F. V., Dornhege, C., and Burgard, W. (2018). Crop Row Detection on Tiny
- Plants With the Pattern Hough Transform. *IEEE Robotics and Automation Letters (RA-L)* 3, 3394–3401
- Wu, X., Aravecchia, S., Lottes, P., Stachniss, C., and Pradalier, C. (2020). Robotic weed control using
  automated weed and crop classification. *Journal of Field Robotics (JFR)* 37, 322–340
- Zenkl, R., Timofte, R., Kirchgessner, N., Roth, L., Hund, A., Van Gool, L., et al. (2022). Outdoor plant
  segmentation with deep learning for high-throughput field phenotyping on a diverse wheat dataset. *Frontiers in Plant Science* 12. doi:10.3389/fpls.2021.774068
- Zhao, L., Zhao, Y., Liu, T., and Deng, H. (2023). A weakly supervised semantic segmentation model of
   maize seedlings and weed images based on scrawl labels. *Sensors* 23. doi:10.3390/s23249846

## **FIGURE CAPTIONS**



**Figure 1** The overview of our pipeline to generate semantic labels for images of crop fields. We use a robotic platform equipped with an RGB camera to collect posed images of the field. Each image gets processed by our automatic labeling method, generating a semantic segmentation of the image to fuse into the semantic map. At each time step, we use the current semantic map to generate the image's semantic label and update the map accordingly.



**Figure 2** Example of a typical UAV mission. The coverage path along which we fuse semantic image labels is depicted in white, the square is the initial pose, and the arrows indicate the direction of movement. The images can overlap, but it is not required. This path maximizes the crop field coverage and is typically used in aerial data collection missions.



**Figure 3** Flowchart of our automatic labeling approach with an example image. At time step t, we take as input the RGB image  $I_t$  recorded from pose  $p_t$  and the set of previously detected lines  $\mathcal{R}_{t-1}$ , depicted in blue boxes. First, we extract the vegetation mask  $I_{t,vm}$  using a graph-segmentation approach (Felzenszwalb and Huttenlocher, 2004). Based on  $I_{t,vm}$ , we compute the connected components to extract plant instances and compute the most prominent crop row via the Hough transform. We propagate the set of previously detected crop rows  $\mathcal{R}_t$  into the current image  $I_t$  to track multiple crop rows. The newly detected crop row in  $I_t$  is added to  $\mathcal{R}_t$ . Then, we label all connected components in  $I_{t,vm}$  as crops that intersect one of the crop row in  $\mathcal{R}_t$  and assign them to the weed class if their distance is above a certain threshold. Vegetation components which are too close to detected crop rows are assigned an "unknown" class that is ignored during network training to minimise labeling errors and thus maximise prediction performance.

#### **Frontiers**



**Figure 4** Given the vegetation mask, we visualize the line detected by the Hough transform (in white). Considering the origin as the bottom left corner, we show the parameters  $r_i$  and  $\theta_i$  defining the detected lines. The vegetation components intersecting the line are thus labeled as crop (green). We can see a weed (red) on the left of the image, since the vegetation component is far away from the detected line.



Ground Truth

Network prediction

Uncertainty

Post-processed prediction



**Figure 5** For the RGB input image on the left, we show the semantic ground truth labels, where crops are represented in green, weeds in red, and the soil in black. Then, we show our network's prediction and we highlights some mistakes using white dotted circles, where weeds are mislabeled as crops. The fourth image shows our network's uncertainty. As expected, the network is mostly uncertain about the boundaries of the plants and about the weeds, we wee that even the weeds labeled as crops in our prediction have high uncertainty. The last image shows our post-processed prediction, after we label as weeds the highly uncertain vegetation components. We can see that this corrects many of the network's errors.



**Figure 6** Qualitative results of our and Lottes et al. methods on PhenoBench (top row) and Onions (bottom row). Soil is black, crops are green, weeds are red, vegetation that we leave unlabeled is white. In the dashed blue circles, we highlight segmentation errors.

#### **TABLES**

**Table 1** Details for the datasets used in the paper: name, reference paper, camera sensor, image resolution and GSD.

Dataset	Reference	Camera	Image Resolution [px]	$\text{GSD}\left[ \tfrac{mm}{px} \right]$
PhenoBench	Weyler et al. (2024)	PhaseOne iXM-100 with a 80 mm RSM prime lens on a gimbal (UAV)	$11664{ imes}8750$	1
Carrots	Bosilj et al. (2020)	Teledyne DALSA Genie Nano deployed on a manually pulled cart (UGV)	$2428\!\times\!1985$	0.4
Onions	Bosilj et al. (2020)	Teledyne DALSA Genie Nano deployed on a manually pulled cart (UGV)	$2149{\times}1986$	0.4
Sugar Beets	Weyler et al. (2022b)	PhaseOne iXM-100 (UAV)	$4320{\times}4100$	1.5

Table 2 List of the hyperparameters of our method, where they are used, and their chosen values.

Hyperparameter	Method	Value		
minimum number of pixels for detection $(\tau_{px})$	Hough line detection	H (i.e. image height)		
width of the line to fit $(l_w)$	Hough line detection	$5\mathrm{px}$		
confidence interval for crop rows ( $\delta$ )	Weed labeling	3		
maximum number of annealing epochs $(T)$	Evidential Deep Learning	25		

**Table 3** Performance of all the baselines on the PhenoBench dataset, Carrots dataset, Onion dataset, and Sugar Beets dataset. The top rows are the general purpose approaches, while the bottom rows are the domain-specific ones. We report the mean IoU, plus the IoU and boundary IoU per class. In **bold** the best results per column.

Dotoset	Approach	IoU [%]			mIoU	Boundary IoU [%]		
Dataset	Approach	soil	crop	weed	· miou	soil	crop	weed
	STEGO	21.4	11.9	0.4	11.2	0.0	1.5	0.0
DhanoBanch	U2Seg	84.6	40.0	2.4	42.3	45.8	11.7	3.4
Thenobelien	Lottes (base)	99.6	44.1	7.6	50.5	0.0	0.0	0.9
	Ours (base)	98.8	80.7	7.2	62.2	86.3	79.1	13.2
	STEGO	28.4	5.1	15.8	16.4	0.0	0.9	0.0
Carrots	U2Seg	80.1	20.4	2.3	34.3	36.2	0.0	19.3
Carrots	Lottes (base)	89.1	15.9	34.0	46.3	0.0	0.0	6.8
	Ours (base)	90.4	12.6	42.7	48.6	84.4	23.6	9.4
	STEGO	26.5	5.1	3.0	11.5	0.0	2.4	0.0
Onion	U2Seg	92.8	0.0	4.3	32.4	24.2	0.0	8.2
Onion	Lottes (base)	89.7	1.4	1.1	30.7	0.0	0.0	1.6
	Ours (base)	95.4	10.7	16.6	40.9	74.2	10.7	16.7
	STEGO	24.9	4.7	1.3	10.3	0.0	1.9	0.0
Sugar Beets	U2Seg	77.9	9.9	6.7	31.5	1.8	2.8	0.0
Sugar Decis	Lottes (base)	98.0	23.6	18.8	46.8	0.0	0.0	1.5
	Ours (base)	97.7	50.6	24.7	57.7	88.7	0.0	1.8

**Table 4** Performance of ERFNet trained on the labels generated by ours and the approach by Lottes et al. We also report the results when we use the uncertainty to post process the semantic predictions. The bottom line shows a fully supervised approach trained on manual labels as upper bound of the performance. Best results per column in **bold**.

Approach		mIoII		
Арргоасн	soil	crop	weed	mioc
Lottes et al. (learned)	99.1	54.6	11.2	55.0
Lottes et al. (+ uncertainty)	99.1	27.2	8.1	44.8
Ours (learned)	99.1	88.8	21.0	69.6
Ours (+ uncertainty)	99.1	88.6	22.7	70.1
Ours (PhenoBench test)	99.5	87.9	24.6	70.7
ERFNet (fully supervised)	98.0	83.4	33.5	71.6

**Table 5** Performance of fully supervised models trained on manually annotated data, and in brackets the difference with respect to the model after fine-tuning. In red if the fine-tuned model performs worse, in blue if it performs better. The gray cells show the performance on the same dataset.

Train		Test	IoU [%]						mIoI	
		Test	soil		crop		weed		mioo	
PhenoBench	(+ Sugar Beets)	PhenoBench	98.0	(-0.4)	83.4	(-11.0)	33.5	(-11.7)	71.6	(-7.7)
		Sugar Beets	93.5	(+0.2)	7.3	(+44.4)	16.8	(+8.2)	39.2	(+17.6)
		Carrots	89.0	(-2.5)	11.1	(+14.9)	47.1	(-11.7)	49.1	(+0.2)
		Onions	82.4	(+5.3)	0.5	(+5.0)	11.3	(-4.4)	31.4	(+2.0)
Sugar Beets	(+ PhenoBench)	PhenoBench	97.6	(-0.1)	67.0	(+9.8)	11.7	(+4.7)	60.2	(+3.4)
		Sugar Beets	98.3	(-4.2)	72.4	(-10.9)	59.2	(-20.5)	76.6	(-11.8)
		Carrots	87.6	(+1.0)	36.1	(+2.1)	24.3	(+10.0)	49.0	(+4.7)
		Onions	86.3	(+1.0)	0.2	(+12.1)	13.2	(+0.7)	33.2	(+4.6)