# Weakly Supervised Minirhizotron Image Segmentation with MIL-CAM

Guohao Yu[1] (iD), Alina Zare[1] (iD), Weihuang Xu[1] (iD), Roser Matamala[2] (iD), Joel Reyes-Cabrera[3] (iD), Felix B. Fritschi[3] (iD), and Thomas E. Juenger[4] (iD)

[1] Department of Electrical and Computer Engineering, University of Florida, Gainesville, FL, USA 32611 `azare@ece.ufl.edu`
[2] Argonne National Laboratory, Lemont, IL, USA 60439
[3] Division of Plant Sciences, University of Missouri, Columbia, MO, USA 65211
[4] Department of Integrative Biology, University of Texas at Austin, Austin, TX, USA 78712

**Abstract.** We present a multiple instance learning class activation map (MIL-CAM) approach for pixel-level minirhizotron image segmentation given weak image-level labels. Minirhizotrons are used to image plant roots *in situ*. Minirhizotron imagery is often composed of soil containing a few long and thin root objects of small diameter. The roots prove to be challenging for existing semantic image segmentation methods to discriminate. In addition to learning from weak labels, our proposed MIL-CAM approach re-weights the root versus soil pixels during analysis for improved performance due to the heavy imbalance between soil and root pixels. The proposed approach outperforms other attention map and multiple instance learning methods for localization of root objects in minirhizotron imagery.

## 1 Introduction

Minirhizotron (MR) imaging plays an important role in plant root studies. It is a widely-used non-destructive root sampling method used to monitor root systems over extended periods of time without repeatedly altering critical soil conditions or root processes [11,4,31,24]. Yet, a significant bottleneck which impacts the value of MR systems is the analysis time needed to process collected imagery. Standard analysis approaches require manual root tracing and labeling of root characteristics. Manually tracing roots collected with MR systems is very tedious, slow, and error prone. Thus, MR image analysis would greatly benefit from automated methods to segment and trace roots. There have been advancements made in this area [33,36,38]. However, the effective automated methods still require a manually-labeled training set. Although these approaches provide a reduction in effort needed over hand-tracing an entire collection of data, the generation of these training sets is still time consuming and labor intensive. In this paper, we propose a weakly supervised MR image segmentation method that relies only on image-level labels.

arXiv:2007.15243v1 [cs.CV] 30 Jul 2020

By relying only on weak image level labels (e.g., this image does/does not contain roots), the time and labor needed to generate a training set is drastically reduced [21,22,42,26,35,30]. It is also much easier and less error prone to identify when an image does or does not contain roots as opposed to correctly labeling every pixel in an image. However, current weakly-supervised methods used to infer pixel-levels labels do not perform as well as semantic image segmentation methods leveraging full annotation [17,25,16,6,3].
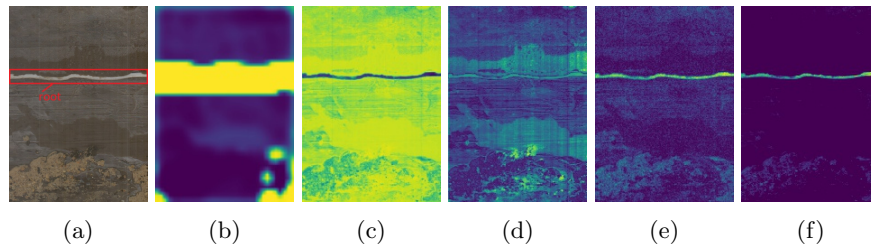


Fig. 1: Example attention maps from different methods of semantic segmentation of an MR image. (a) Original MR Image. (b) CAM Result (c) Grad-CAM Result (d) Grad-CAM++ Result (e) SMOOTHGRAD Result (f) Result of proposed method, MIL-CAM.

Attention or class activation maps have been widely used to infer pixel-level labels from training data with weak image-level labels [42,27,5,28]. However, existing attention map approaches have been found to be inaccurate in identifying and delineating roots in MR imagery. For example, CAM (the class activation maps) approach [42] overestimates the size of the roots as shown in Fig.1b. The Gradient-weighted Class Activation Mapping (Grad-CAM) [27] approach incorrectly identifies the background soil as the root target as shown in Fig.1c. Grad-CAM++ [5] and SMOOTHGRAD [28] shown in Fig.1d and Fig.1e, respectively, result in maps with poor contrast between roots and soil and, thus, many false alarms.

In this paper, we propose the multiple instance learning CAM (MIL-CAM) approach to address root segmentation in MR imagery given weak image-level labels. MIL-CAM is outlined in Section 3. In Section 4, we compare MIL-CAM approach results to existing approaches with both weak- and full-annotation on an MR dataset collected from switchgrass.

## 2    Related Work

### 2.1    Attention Maps for Semantic Segmentation

CAM [42] is one of the earliest methods showing that attention maps can localize the discriminative image components for classification. CAM uses a network
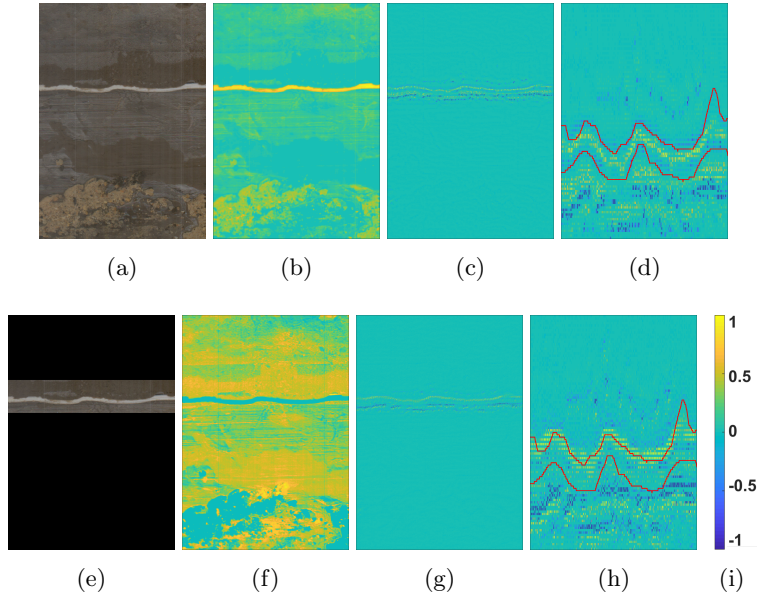
Fig. 2: Examples of the gradients with respect to the image classification score of the target root class using various individual feature maps. (a) Original MR image. (b) Feature map of Fig.2a which highlights the root area. (c) Gradients of feature map in Fig.2b with respect to the root class score. (d) Gradients in Fig.2c around the masked root regions shown in Fig.2e. The red lines indicate the boundary of root object. (e) Cropped area from full image as mask in Fig.2d and Fig.2h. (f) Feature map of Fig.2a which highlights soil. (g) Gradients of feature map in Fig.2g. (h) Gradients in Fig.2g around root object in Fig.2e. The red lines indicate the boundary of the root object. Fig.2d and Fig.2h are rescaled to match the size of other subfigure. (i) Colorbar used in for images in Fig. 2.

structure consisting of a block of fully convolutional layers followed by a global average pooling layer and a single fully connected layer. The block of fully convolutional layers extract image features. These extracted features are combined linearly using the weights of the final fully connected layer to define the CAM. However, the CAM from this approach often has a low resolution, making it challenging to infer pixel-level labels precisely for semantic segmentation. Grad-CAM [27] is an extension of CAM which can estimate higher resolution attention maps by using features from any convolutional layer in the network. Specifically, Grad-CAM estimates the weights used for combining features as the average of gradients of the image classification score with respect to each value in the corresponding feature maps. Following the introduction of Grad-CAM, many methods were proposed to attempt to improve the quality of attention maps generated. Grad-CAM++ [5] takes a weighted average of only the positive gradients of the

image classification score with respect to each feature map. SMOOTHGRAD [28] averages over several Grad-CAM estimated attention maps with added zero-mean Gaussian noise with the aim of reducing sensitivity to feature map noise. Smooth Grad-CAM++ [18] mimics SMOOTHGRAD but applies the approach to Grad-CAM++ estimated attention maps. Score-CAM [32] attempts to improve Grad-CAM by weighting feature maps based on a metric which measures the increase of confidence for a class associated with the inclusion of each feature map. In application, all of these methods have been found to be either imprecise or sensitive to imbalanced data sets. Specifically in our application, soil pixels having complex gradients (i.e., both positive and negative gradients) which has a huge impact on the weights. Consider the example shown in Fig.2. Gradients across the feature maps have differing signs as shown in Fig.2d and Fig.2h and, thus, when averaged over the map may cancel each other out. Given this, the standard Grad-CAM approach is ineffective since the average of the gradients over the feature map is used to compute the attention map.

## 2.2   Weakly Supervised Learning

Weakly supervised and multiple instance learning (MIL) algorithms for image segmentation do not require precise pixel-level labels. Under MIL, a set of samples (e.g., an image) is labeled as either "positive" or "negative." Positively labeled images are assumed to have at least one pixel corresponding to the target class (i.e., in our case, roots). The number of target pixels in positively labeled images are unknown. Negatively labeled images are composed of only non-target class (i.e., soil) pixels. Often, MIL approaches iteratively estimate the likelihood each pixel is a target and, using these values, update classifier parameters (and, then, subsequently update likelihood values again) [19,7]. Similarly, the pixels with the lowest target likelihood in each positively labeled image is also commonly assumed to be from the non-target class [7,8]. In contrast to methods that select likely target and non-target pixels, some methods have been proposed which consider all pixels in an image as equally contributing to the image-level label [42]. The Log-Sum-Exp (LSE) algorithm uses a hyper-parameter which trades off between selecting a single pixel as the target representative and considering all pixels in an image as target with equal contribution [22]. Global weighted rank pooling (GWRP) is another way to generalize number of pixels identified as targets [12]. In all of these approaches, it is difficult to select a fixed number of pixels to identify as targets representatives.

One reason that identifying the number of target pixels is challenging is that the size of the target class objects vary across images (e.g., some images contain only very few thin, fine roots whereas others are filled with roots of varying diameter). To alleviate this challenge, some approaches identify target pixels by adapting a threshold [35,10,1,14]. In [35], pixels with target class scores larger than a predefined threshold are labeled as targets and pixels with low saliency values are considered background. However, in this approach pixels are often unassigned to either target or background classes, pixels may be assigned to multiple target labels, or pixels may be assigned to background despite being

surrounded in their neighborhood by target pixels. In [34], pixels in all of these cases are ignored. The approach outlined in [10] attempts to deal with these ignored pixels using deep seeded region growing (DSRG). DSRG proposes to propagate labels from labeled pixels to unlabeled pixels. The method presented in [14] extends the DSRG approach [10] by thresholding aggregated localization maps to improve delineation of target regions and adapts their algorithm to accommodate semi-supervised segmentation. Following an initial segmentation, some approaches apply post-processing steps to smooth and improve segmentation labels. These include conditional random fields (CRF) and the GWRP approach [13,1,12].

### 2.3  MR Image Segmentation

Several methods have been developed for automated minirhizotron image segmentation [40,9,41,23]. Currently, supervised deep learning approaches are the methods that are achieving the state-of-art results in MR image segmentation [36,33,37,29]. Yet, deep learning methods require a large collection of precisely traced root images for training the networks. A small number of approaches have been investigated for weakly supervised MR image segmentation [38]. Yu, et al. [38] studied the application of three MIL algorithms: multiple instance adaptive cosine coherence estimator (MI-ACE) [39], multiple instance support vector machine (miSVM) [2], and multiple instance learning with randomized trees (MIForests) [15] for application to MR imagery. These methods, however, did not do feature learning and, so, the authors manually identified color features to be used during segmentation.

## 3  MIL-CAM Methodology

Semantic segmentation from weak labels using MIL-CAM is achieved in two training stages. The first stage, outlined in Alg. 1, estimates the set of parameters needed to compute an attention map $\mathbf{S}^c \in \mathbb{R}^{M \times N}$ for a class $c$ where $M$ and $N$ are the numbers of rows and columns of the input image, respectively. The attention map is estimated using the softmax output of a weighted linear combination feature maps extracted from the various layers of a trained CNN as described in Eq. 1,

$$\mathbf{S}^c = \frac{exp(\sum_j w_j^c y(\mathbf{F}_j) + b^c)}{\sum_q exp(\sum_j w_j^q y(\mathbf{F}_j) + b^q)}. \tag{1}$$

where $q$ is an index over all output classes, $w_j^q \in \mathbb{R}$ is the weight estimated for class $q$ and feature map $\mathbf{F}_j \in \mathbb{R}^{A_j \times B_j}$, $A_j$ and $B_j$ are the number of rows and columns in the $j^{th}$ feature map, $b^q \in \mathbb{R}$ is an estimated bias term, and $y(\cdot)$ is an interpolation function to scale an input to the size of $M \times N$.

Once attention maps are obtained using Alg. 1, a segmentation network is then trained as outlined in Alg. 2. After training, the segmentation network maps input test imagery to get pixel level segmentation outputs.

### 3.1   Attention Map Estimation

MIL-CAM estimates the set of parameters needed to obtain attention maps and compute Eq. 1 using the combination of three key components: (a) a pixel-level feature extraction component; (b) a pixel sampling component used to form a bag for each image for MIL analysis; and (c) a linear model that performs the MIL-based segmentation. The sampled pixels with features extracted from the image classification network are used to train the linear model. The approach is illustrated in Fig. 3 and outlined in the following sub-sections.
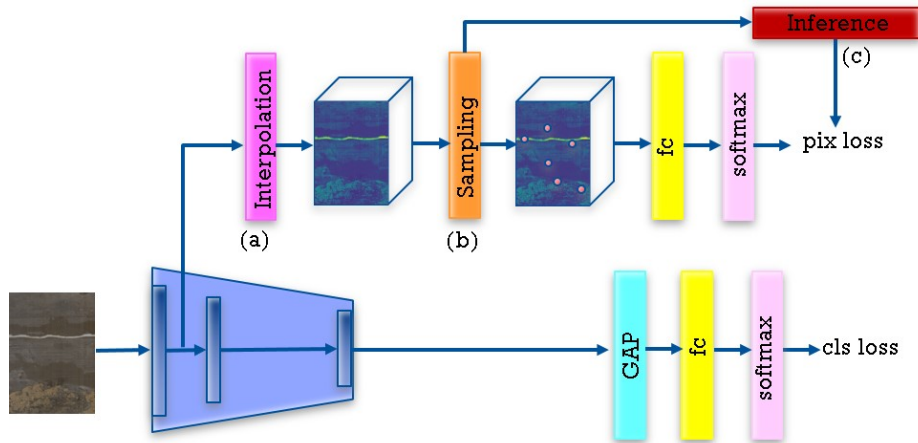


Fig. 3: Architecture of MIL-CAM. GAP represents a global average pooling layer and fc represents a fully connected layer. cls loss represents the loss for image classification into positive (i.e., containing roots) or negative (i.e., does not contain roots). pix loss represents the loss for pixel level classification into root vs. soil.

**Feature Extraction and Interpolation** An image-level CNN classification network is first trained to extract coarse feature maps for each image. The training data set, $\{(\mathbf{I}_1, l_1), ...(\mathbf{I}_k, l_k), ...(\mathbf{I}_K, l_K)\}$, consists of $K$ images where each image $\mathbf{I}_k \in \mathbb{R}^{3 \times M \times N}$ is paired with image label $l_k \in \{0, 1\}$ where 0 represents a negative image (i.e., does not contain roots) and 1 represents a positive image (i.e., contains roots). Using this training data an image-level classification network is trained by optimizing the cross-entropy loss as shown in Eq. 2,

$$\min_{\boldsymbol{\theta}_0} \sum_{k=1}^{K} L_{cls-loss}(\mathbf{I}_k; \boldsymbol{\theta}_0, \mathbf{l}_k) = \frac{-1}{K} \sum_{k=1}^{K} \sum_{q=1}^{2} l_{kq} \log f_q(\mathbf{I}_k; \boldsymbol{\theta}_0) \qquad (2)$$

where $\mathbf{l}_k = [l_{k1}, l_{k2}]$ is the one-hot encoded label of the image $\mathbf{I}_k$ and $f_q(\mathbf{I}_k; \boldsymbol{\theta}_0)$ is the $q^{th}$ element of the softmax output layer of the the image classification network defined by parameters $\boldsymbol{\theta}_0$. The assumption is, provided an effective image-level classification network can be trained, that the network is extracting features that are useful for the semantic segmentation problem and these useful features are encoded in the CNN feature maps. Once the classification network is trained, then the coarse CNN feature maps are upsampled using bilinear interpolation to match the size of the input image. Each pixel is represented by the corresponding feature vector obtained from the collection of upsampled feature maps.

---

**Algorithm 1:** Estimating Weights and Biases for Attention Maps

**Data:** $(\mathbf{I}, \mathbf{l}) = \{(\mathbf{I}_1, l_1), ...(\mathbf{I}_k, l_k), ...(\mathbf{I}_K, l_K)\}$
Train the image classification network with $(\mathbf{I}, \mathbf{l})$ ;
Extract feature maps $\mathbf{F}$ from image classification network for $(\mathbf{I}, \mathbf{l})$ ;
Interpolate the CNN feature maps $y(\mathbf{F})$ for $(\mathbf{I}, \mathbf{l})$;
Sample instances and construct bags, $\{(\mathbf{B}_1, l_1), ...(\mathbf{B}_k, l_k), ...(\mathbf{B}_K, l_K)\}$;
Initialize each instance label with the label of its corresponding bag;
**repeat**
  Update $\mathbf{w}, \mathbf{b}$ by optimizing Eq. 3 with stochastic gradient descent for one
    epoch using the instances and updated labels $(\mathbf{x}_k^n, \mathbf{l}_k^n)$;
  Compute $p_k^n = g(\mathbf{x}_k^n; \mathbf{w}, \mathbf{b})$ for each instance;
  **for** Every positive bag $(\mathbf{B}_k, l_k = 1)$ **do**
    $p_t = \text{Otsu's}(\{p_k^1, ...p_k^{N_k}\})$;
    If $p_k^n \geq p_t$, then set $\mathbf{l}_k^n$ as target, else set $\mathbf{l}_k^n$ as non-target;
  **end**
  **for** Every negative bag $(\mathbf{B}_k, l_k = 0)$ **do**
    set $\mathbf{l}_k^n$ as non-target;
  **end**
**until** Fixed number of epochs completed;
**return** $\boldsymbol{\theta}_0, \mathbf{w}, \mathbf{b}$ from epoch with smallest loss

---

**Instance Sampling** In order to address some of the imbalance in the data set (i.e., there are many more soil pixels than root pixels), a sampling approach is used to identify representative pixels from each image. The green band of the RGB minirhizotron image is used for instance sampling. The approach draws a single pixel to represent the set of pixels from each possible 8-bit value from the green band in the image. In other words, a 256 bin histogram is built using the values of the green band of the MR imagery. For each non-empty bin, a uniform random draw is used to identify a representative pixel for that green-level. In our application, we found this to be an effective approach to re-balance root-vs-nonroot pixels in positively labeled imagery (given that pixel level labels are unavailable). The sampled pixels are organized into a set of

bags, $\{(\mathbf{B}_1, l_1), ...(\mathbf{B}_k, l_k), ...(\mathbf{B}_K, l_K)\}$. Each bag, $\mathbf{B}_k = \left\{\mathbf{x}_k^1, \mathbf{x}_k^2, \ldots, \mathbf{x}_k^{N_k}\right\}$, corresponds to one image $\mathbf{I}_k$ with image label $l_K$ and is composed of $N_k$ instances. The instance $\mathbf{x}_k^n \in \mathbb{R}^J$ is the feature vector for the $n^{th}$ instance in the $k^{th}$ bag where $J$ is the number of feature maps used to construct the feature vectors.

---

**Algorithm 2:** Weakly Supervised Image Segmentation

**Data:** $\{(\mathbf{I}_1, l_1), ...(\mathbf{I}_k, l_k), ...(\mathbf{I}_K, l_K)\}$
**Parameter:** $s_t$
**for** Every positive image **do**
   Compute the score-map $\mathbf{S}_k^c$ from MIL-CAM using $(\boldsymbol{\theta}_0, \mathbf{w}, \mathbf{b})$ in Eq.1 ;
   Estimate a threshold, $o_t = $ Otsu's $(\mathbf{S}_k^c)$;
   If $\mathbf{S}_k^c(m, n) \geq o_t$, then set $\mathbf{l}_k^0(m, n)$ as target, else set $\mathbf{l}_k^0(m, n)$ as non-target;
**end**
**for** Every negative image $(\mathbf{I}_k, l_k = 0)$ **do**
   set $\mathbf{l}_k^0(m, n)$ as non-target;
**end**
Update parameters $\boldsymbol{\theta}_1$ for data set with pixel labels $\mathbf{l}_k^0(m, n)$ for a fixed number of epochs;
**repeat**
   Compute score-map for each image using the U-Net with updated
     parameters;
   **if** A positive image $(\mathbf{I}_k, l_k = 1)$ **then**
      If $\mathbf{P}_k(m, n) \geq s_t$, then set $\mathbf{l}_k(m, n)$ as target, else $\mathbf{l}_k(m, n)$ as non-target;
      If every pixel $\mathbf{P}_k(m, n) < s_t$ , then $\mathbf{l}_k(m, n) = \mathbf{l}_k^0(m, n)$;
   **else if** A negative image $(\mathbf{I}_k, l_k = 0)$ **then**
      set $\mathbf{l}_k(m, n)$ as non-target;
   Update parameters $\boldsymbol{\theta}_1$ for dataset with pixel labels $\mathbf{l}_k(m, n)$ ;
**until** Fixed number of epochs completed;
**return** Segmentation network parameters $\boldsymbol{\theta}_1$

---

**Estimated Weights and Biases** After instance sampling, the weights and biases used to compute the attention maps as defined in Eq. 1 are estimated by optimizing the cross-entropy loss shown in Eq. 3 given the MIL constraints that for each positive bag, at least one instance must be labeled as root and all instances in every negative bag are labeled as non-root,

$$\min_{\mathbf{l}_k^n} \min_{\mathbf{w}, \mathbf{b}} \sum_{\mathbf{x}_k^n} L_{pix-loss}(\mathbf{x}_k^n; \mathbf{w}, \mathbf{b}, \mathbf{l}_k^n) = \frac{-1}{\sum_k N_k} \sum_{\mathbf{x}_k^n} \sum_q l_{kq}^n \log g_q(\mathbf{x}_k^n; \mathbf{w}, \mathbf{b}) \quad (3)$$

where $\mathbf{l}_k^n$ is the one-hot encoded label of the instance $\mathbf{x}_k^n$ and $g_q(\mathbf{x}_k^n; \mathbf{w}, \mathbf{b})$ is the $q^{th}$ element of the softmax output of the MIL-CAM with parameters $(\mathbf{w}, \mathbf{b})$. The loss is updated iteratively as outlined in Alg. 1. During the initial epoch, each instance is labeled the same label as its bag. In all subsequent epochs, the

probability that an instance belongs to the target class, $p_k^n = g(\mathbf{x}_k^n; \mathbf{w}, \mathbf{b})$, is predicted by the linear model trained from the previous epoch. Then for each positive bag, a threshold $p_t$ is computed using Otsu's threshold [20] and all instances greater than the threshold are labeled as target whereas all others are labeled as non-target. For negative bags, all instances are labeled as non-target.

### 3.2 Training the Image Segmentation Network

Once MIL-CAM attention maps can be estimated, an image segmentation network is trained as outlined in Alg.2. First, target class attention maps for positively labeled images are estimated and thresholded using Otsu's threshold to obtain a label for each pixel. All pixels in negatively labeled images are given a non-target label. These labels are used to estimate the parameters for the U-Net [25] architecture illustrated in lower branch of Fig. 4. After initially training the U-Net with labels obtained from the attetnion maps, the U-Net is iteratively fine-tuned. A score-map, $\mathbf{P}_k \in \mathbb{R}^{M \times N}$, is computed using the soft-max output of the U-Net. The score-map of positively-labeled is thresholded using a fixed (large) threshold parameter, $s_t$, to obtain updated pixel level labels which highlight more likely positive samples. The updated labels are iteratively used to fine-tune the parameters of the U-Net.
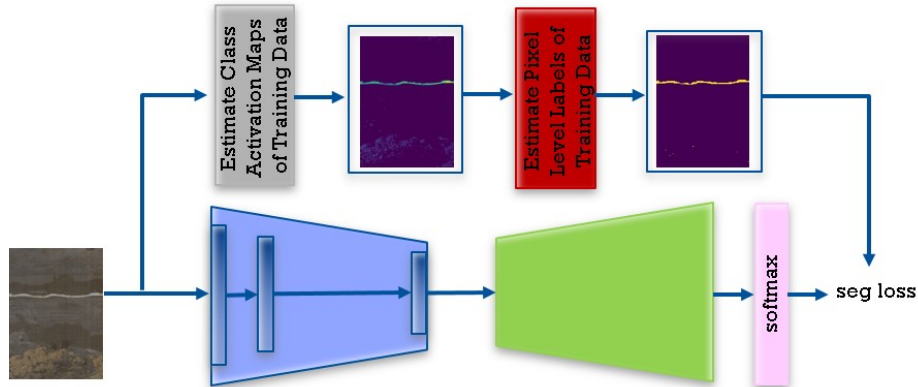


Fig. 4: Architecture of segmentation U-Net with MIL training branch. The bottom branch is the U-Net. The top branch is used to infer label of training data.

## 4 Experiments

### 4.1 Data Description

For our experiments, we used a switchgrass (*Panicum virgatum* L.) MR imagery dataset consisting of 561 training images with image-level labels and 30 test and

validation images with pixel-level labels. Each image was $2160 \times 2550$ in size and was divided into sub-images of size $720 \times 510$. 500 sub-images containing roots and 500 sub-images containing only soil were randomly selected as training data for estimating attention map parameters. 1500 root sub-images and 1500 soil sub-images were randomly selected as training data for the U-Net segmentation network. The 30 images with pixel-level labels were randomly divided into 10 validation images and 20 test images.

## 4.2   Architecture

Our experiments use U-Net [25] with layer depth of 5 as backbone for MR image segmentation. The feature extraction network used to estimate attention map parameters was a 2-class convolutional neural network with the encoder of the U-Net, followed by a global average pooling layer and a fully connected layer. We extract $1024 \times 46 \times 33$ feature maps and vectorize the feature maps to classify each image into 2 classes with a fully connected layer. The feature extraction net is trained using SGD at a learning rate of 0.0001 and momentum of 0.8 in the online mode to minimize the cross entropy loss. The MIL-CAM attention map module extracts a 64-dimensional feature for each sampled instance from the fourth layer of the encoder of the feature extraction network. Then, classifies each sampled instance into one of two classes using a fully connected layer. The MIL-CAM attention map module is trained using SGD at a learning rate of 0.001 and momentum of 0.5 in the online mode to minimize the cross entropy loss.

The image segmentation network was a U-Net of depth 5 and a MIL training branch. The MIL training branch extracts $64 \times 720 \times 510$ features from the first layer of the encoder of the feature extraction network and compute a $720 \times 510$ score-map of target class for each training image. The threshold parameter $s_t$ was set to 0.9 to estimate pixel label from the score-map. The U-Net was first initialized for 10 epochs using Adam at learning rate of 0.0001 in the online mode to minimize the cross entropy loss where the root class was weighted by 50 using the labels produced by the attention maps. Then, during iterative fine-tuning, the network parameters were also updated using Adam with learning rate of 0.0001 in the online mode to minimize the cross entropy loss with the root class having an additional weight of 50. The weight on root class addressed the imbalance issue between root class and soil class.

## 4.3   Experiments: MIL-CAM Attention Maps

The attention maps of MIL-CAM were first qualitatively compared with attention maps of other methods as shown in Fig.5. As can be seen, MIL-CAM results shown in Fig.5f more accurately indicate root locations as compared to the attention maps produced by CAM in Fig.5b. This difference in performance is largely due to the fact that CAM requires interpolating a low resolution attention map to the size of the input image resulting in blurred, oversized detection regions. Grad-CAM in Fig.5c. fails to correctly identify roots and, instead, highlights

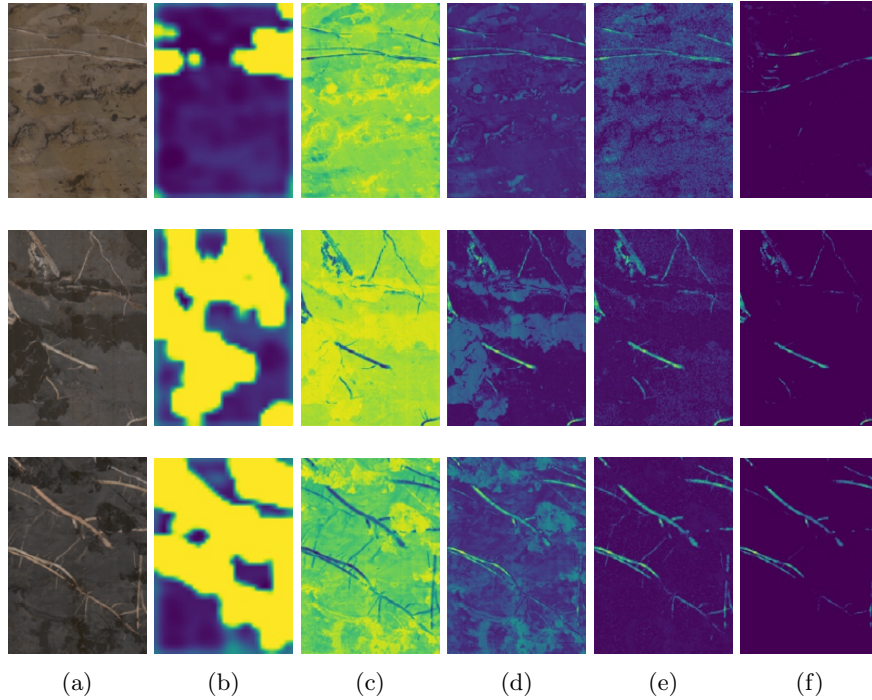|  (a)  |  (b)  |  (c)  |  (d)  |  (e)  |  (f)  |

Fig. 5: Attention maps of different methods. (a) Original Image. (b) Result of CAM. (c) Result of Grad-CAM. (d) Result of Grad-CAM++. (e) Result of SMOOTHGRAD. (f) Result of MIL-CAM.

soil. Furthermore, MIL-CAM produced attention maps with higher contrast between root pixels and background than those Grad-CAM ++ in Fig.5d and SMOOTHGRAD in Fig.5e.

Fig.6 compares attention maps from a selection of approaches after thresholding with Otsu's threshold. Table 1 lists the average and standard devation for precision, recall and F1 score of three training runs of the various approaches to compare the quality these thresholded results. The proposed MIL-CAM method has a significantly higher F1 score among all those compared. The precision of MIL-CAM is an order of magnitude better than the comparison methods without a significant loss in recall as compared with the gains of precision. Although other methods except Grad-CAM have a better recall, the low precision scores of these methods indicate a large amount of background pixels are mislabeled as root pixels. This can be visualized in Fig.6.

### 4.4 Experiments: Semantic Segmentation

We also compared the performance of our final MIL segmentation network (i.e., MIL-CAM Th in the table) against other MIL methods (MI-ACE[38],
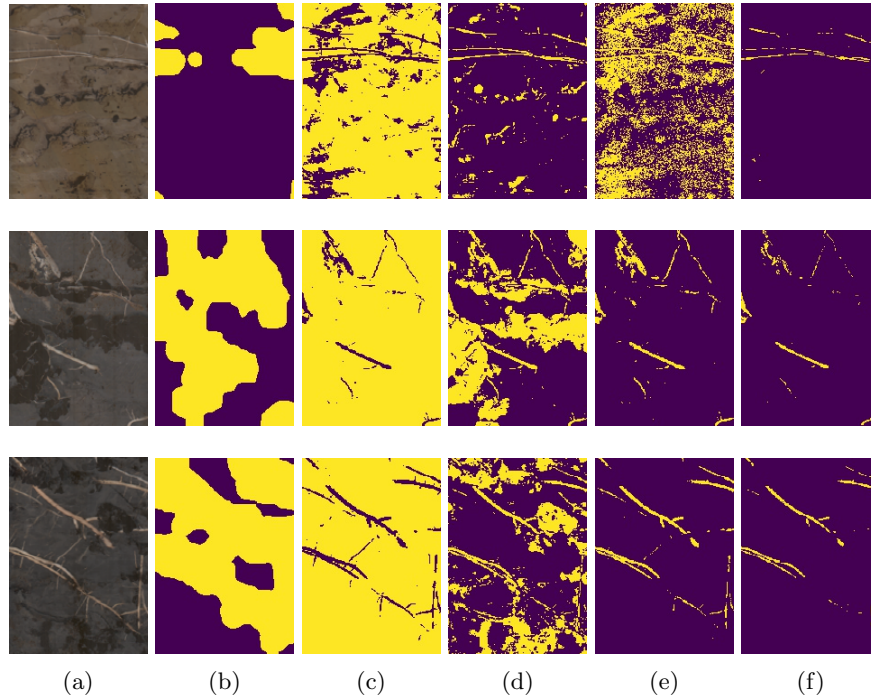
Fig. 6: Thresholded attention maps. (a) Original Image. (b) Result of CAM. (c) Result of Grad-CAM. (d) Result of Grad-CAM++. (e) Result of SMOOTH-GRAD. (f) Result of MIL-CAM.

miSVM[38], and MIForest[38]). The average and standard deviation of three runs of the precision, recall, F1 score and mIoU were compared at false positive rate (FPR) is 0.03 in Table 2. Our proposed approach outperformed all other MIL methods. The proposed MIL-CAM Th method (i.e., the thresholded MIL-CAM result) achieved recall= 0.878. The recall of MIL-CAM Th was 10% better than miSVM which was the second best. MIL-CAM Th also had the best precision of all MIL methods.

The segmentation results of the proposed MIL-CAM approach when taking the argmax of the softmax outputs (i.e., argmax MIL-CAM in the table) are shown in the third column in Fig. 7c. The long roots are a challenging problem. Although our proposed method detects most of the root pixels, it expands the boundary of some roots. This expansion results in high recall (0.859) but low precision (0.186) as shown in table 2. To mitigate this, we also applied a conditional random field (CRF) [13] postprocessing to the segmentation results of our approach. The default parameters of the CRF were used as 0.7 for the certainty of the label, 3 for the parameter of the smoothness kernel, 80 for the spatial parameter of the appearance kernel, 13 for the color parameter of the appearance

Table 1: Compare results of thresholded attention maps

| Method | Precision | Recall | F1 score | mIoU |
|---|---|---|---|---|
| CAM | $0.045 \pm 0.0053$ | $\mathbf{0.931 \pm 0.0459}$ | $0.085 \pm 0.0098$ | $0.045 \pm 0.0053$ |
| Grad-CAM | $0.003 \pm 0.0012$ | $0.229 \pm 0.0939$ | $0.006 \pm 0.0024$ | $0.003 \pm 0.0012$ |
| Grad-CAM++ | $0.015 \pm 0.0084$ | $0.550 \pm 0.1951$ | $0.030 \pm 0.0159$ | $0.015 \pm 0.0083$ |
| SMOOTHGRAD | $0.033 \pm 0.0028$ | $0.782 \pm 0.0191$ | $0.064 \pm 0.0052$ | $0.033 \pm 0.0028$ |
| MIL-CAM | $\mathbf{0.248 \pm 0.1870}$ | $0.536 \pm 0.1450$ | $\mathbf{0.289 \pm 0.1814}$ | $\mathbf{0.177 \pm 0.1190}$ |

Table 2: Comparison of image segmentation results. All comparison methods use weak image level labels except the U-Net approach from [36]. MIL-CAM Th is the result found after thresholding the U-Net softmax outputs corresponding to the target class at FPR = 0.03; argmax MIL-CAM is the result when taking the argmax of U-Net softmax outputs; and MIL-CAM + CRF method is the result when the argmax MIL-CAM result is postprocessed with a CRF.

| Method | Label | Precision | Recall | F1 score | mIoU |
|---|---|---|---|---|---|
| U-Net [36] | pixel | 0.307 | $\mathbf{0.913}$ | 0.459 | 0.298 |
| MI-ACE[38] | image | $0.130 \pm 0.0010$ | $0.775 \pm 0.0067$ | $0.223 \pm 0.0017$ | $0.125 \pm 0.0011$ |
| miSVM[38] | image | $0.134 \pm 0.0015$ | $0.798 \pm 0.0104$ | $0.229 \pm 0.0026$ | $0.129 \pm 0.0017$ |
| MIForests[38] | image | $0.101 \pm 0.0104$ | $0.582 \pm 0.0664$ | $0.172 \pm 0.0180$ | $0.094 \pm 0.0108$ |
| MIL-CAM Th | image | $0.145 \pm 0.0050$ | $0.878 \pm 0.0341$ | $0.249 \pm 0.0088$ | $0.142 \pm 0.0057$ |
| argmax MIL-CAM | image | $0.186 \pm 0.0278$ | $0.859 \pm 0.0423$ | $0.304 \pm 0.0364$ | $0.180 \pm 0.0251$ |
| MIL-CAM + CRF | image | $\mathbf{0.667 \pm 0.0257}$ | $0.692 \pm 0.0267$ | $\mathbf{0.678 \pm 0.0058}$ | $\mathbf{0.513 \pm 0.0066}$ |

kernel and 2 inference steps were run. Segmentation results after CRF postprocessing are shown in the fourth column in in Fig. 7c. Postprocessing improved the precision of results from 0.186 to 0.667, and the mean Intersection-Over-Union (mIoU) from 0.180 to 0.513 as shown in Table 2. The only approach with that outperformed the proposed MIL-CAM with CRF postprocessing on any metric was the U-Net method outlined in [36]. However, this U-Net was pretrained using a large dataset consisting of 17567 MR images with full pixel-level annotation and, thus, did not have to overcome the weak label challenge.

## 5   Conclusion

In this work, we proposed MIL-CAM for weakly supervised MR image segmentation. The proposed MIL-CAM approach outperformed a variety of comparison attention map approaches as well as a variety of MIL segmentation methods, particularly when incorporating a CRF post-processing.
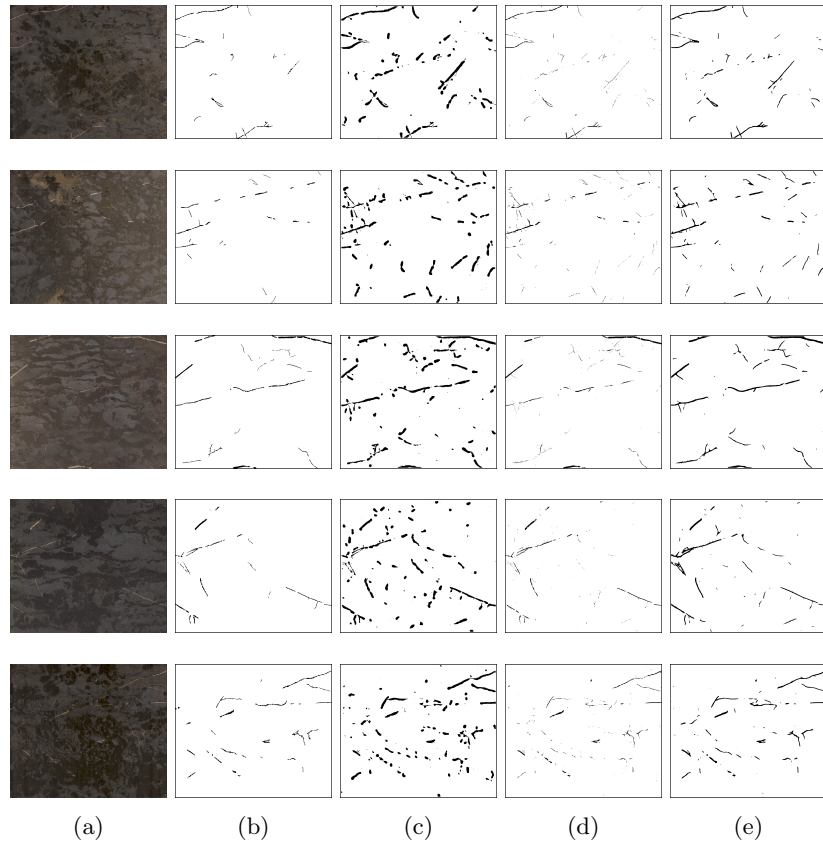
Fig. 7: Qualitative examples of root segmentation results with different method. (a) Original image. (b) groundtruth (GT). (c) Result of argmax MIL-CAM (d) Result of argmax MIL-CAM + CRF. (e) Result of U-Net.

## Acknowledgements

# References

1. Ahn, J., Cho, S., Kwak, S.: Weakly supervised learning of instance segmentation with inter-pixel relations. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2209–2218 (2019)
2. Andrews, S., Tsochantaridis, I., Hofmann, T.: Support vector machines for multiple-instance learning. In: Advances in neural information processing systems. pp. 577–584 (2003)
3. Badrinarayanan, V., Kendall, A., Cipolla, R.: Segnet: A deep convolutional encoder-decoder architecture for image segmentation. IEEE transactions on pattern analysis and machine intelligence **39**(12), 2481–2495 (2017)
4. Bates, G.: A device for the observation of root growth in the soil. Nature **139**(3527), 966–967 (1937)
5. Chattopadhay, A., Sarkar, A., Howlader, P., Balasubramanian, V.N.: Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. In: 2018 IEEE Winter Conference on Applications of Computer Vision (WACV). pp. 839–847. IEEE (2018)
6. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. IEEE transactions on pattern analysis and machine intelligence **40**(4), 834–848 (2017)
7. Durand, T., Mordan, T., Thome, N., Cord, M.: Wildcat: Weakly supervised learning of deep convnets for image classification, pointwise localization and segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 642–651 (2017)
8. Durand, T., Thome, N., Cord, M.: Weldon: Weakly supervised learning of deep convolutional neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4743–4752 (2016)
9. Heidari, M., et al.: A new method for root detection in minirhizotron images: Hypothesis testing based on entropy-based geometric level set decision. International Journal of Engineering **27**(1), 91–100 (2014)
10. Huang, Z., Wang, X., Wang, J., Liu, W., Wang, J.: Weakly-supervised semantic segmentation network with deep seeded region growing. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 7014–7023 (2018)
11. Johnson, M., Tingey, D., Phillips, D., Storm, M.: Advancing fine root research with minirhizotrons. Environmental and Experimental Botany **45**(3), 263–289 (2001)
12. Kolesnikov, A., Lampert, C.H.: Seed, expand and constrain: Three principles for weakly-supervised image segmentation. In: European Conference on Computer Vision. pp. 695–711. Springer (2016)
13. Krähenbühl, P., Koltun, V.: Efficient inference in fully connected crfs with gaussian edge potentials. In: Advances in neural information processing systems. pp. 109–117 (2011)
14. Lee, J., Kim, E., Lee, S., Lee, J., Yoon, S.: Ficklenet: Weakly and semi-supervised semantic image segmentation using stochastic inference. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 5267–5276 (2019)
15. Leistner, C., Saffari, A., Bischof, H.: Miforests: Multiple-instance learning with randomized trees. In: European Conference on Computer Vision. pp. 29–42. Springer (2010)

16. Lin, G., Shen, C., Van Den Hengel, A., Reid, I.: Efficient piecewise training of deep structured models for semantic segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3194–3203 (2016)
17. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3431–3440 (2015)
18. Omeiza, D., Speakman, S., Cintas, C., Weldermariam, K.: Smooth grad-cam++: An enhanced inference level visualization technique for deep convolutional neural network models. arXiv preprint arXiv:1908.01224 (2019)
19. Oquab, M., Bottou, L., Laptev, I., Sivic, J.: Is object localization for free?-weakly-supervised learning with convolutional neural networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 685–694 (2015)
20. Otsu, N.: A threshold selection method from gray-level histograms. IEEE transactions on systems, man, and cybernetics **9**(1), 62–66 (1979)
21. Papandreou, G., Chen, L.C., Murphy, K.P., Yuille, A.L.: Weakly-and semi-supervised learning of a deep convolutional network for semantic image segmentation. In: Proceedings of the IEEE international conference on computer vision. pp. 1742–1750 (2015)
22. Pinheiro, P.O., Collobert, R.: From image-level to pixel-level labeling with convolutional networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1713–1721 (2015)
23. RAHMANZADEH, B.H., Shojaedini, S.: Novel automated method for minirhizotron image analysis: Root detection using curvelet transform. INTERNATIONAL JOURNAL OF ENGINEERING (2016)
24. Rewald, B., Ephrath, J.E.: Minirhizotron techniques. Plant roots: The hidden half pp. 1–15 (2013)
25. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: International Conference on Medical image computing and computer-assisted intervention. pp. 234–241. Springer (2015)
26. Roy, A., Todorovic, S.: Combining bottom-up, top-down, and smoothness cues for weakly supervised image segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3529–3538 (2017)
27. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-cam: Visual explanations from deep networks via gradient-based localization. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 618–626 (2017)
28. Smilkov, D., Thorat, N., Kim, B., Viégas, F., Wattenberg, M.: Smoothgrad: removing noise by adding noise. arXiv preprint arXiv:1706.03825 (2017)
29. Smith, A.G., Petersen, J., Selvan, R., Rasmussen, C.R.: Segmentation of roots in soil with u-net. Plant Methods **16**(1), 1–15 (2020)
30. Tang, M., Djelouah, A., Perazzi, F., Boykov, Y., Schroers, C.: Normalized cut loss for weakly-supervised cnn segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1818–1827 (2018)
31. Waddington, J.: Observation of plant roots in situ. Canadian Journal of Botany **49**(10), 1850–1852 (1971)
32. Wang, H., Du, M., Yang, F., Zhang, Z.: Score-cam: Improved visual explanations via score-weighted class activation mapping. arXiv preprint arXiv:1910.01279 (2019)
33. Wang, T., Rostamza, M., Song, Z., Wang, L., McNickle, G., Iyer-Pascuzzi, A.S., Qiu, Z., Jin, J.: Segroot: A high throughput segmentation method for root image analysis. Computers and Electronics in Agriculture **162**, 845–854 (2019)

34. Wei, Y., Feng, J., Liang, X., Cheng, M.M., Zhao, Y., Yan, S.: Object region mining with adversarial erasing: A simple classification to semantic segmentation approach. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1568–1576 (2017)

35. Wei, Y., Xiao, H., Shi, H., Jie, Z., Feng, J., Huang, T.S.: Revisiting dilated convolution: A simple approach for weakly-and semi-supervised semantic segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 7268–7277 (2018)

36. Xu, W., Yu, G., Zare, A., Zurweller, B., Rowland, D., Reyes-Cabrera, J., Fritschi, F.B., Matamala, R., Juenger, T.E.: Overcoming small minirhizotron datasets using transfer learning. Computers and Electronics in Agriculture **175** (2020). https://doi.org/https://doi.org/10.1016/j.compag.2020.105466

37. Yasrab, R., Atkinson, J.A., Wells, D.M., French, A.P., Pridmore, T.P., Pound, M.P.: Rootnav 2.0: Deep learning for automatic navigation of complex plant root architectures. GigaScience **8**(11), giz123 (2019)

38. Yu, G., Zare, A., Sheng, H., Matamala, R., Reyes-Cabrera, J., Fritschi, F.B., Juenger, T.E.: Root identification in minirhizotron imagery with multiple instance learning. Machine Vision and Applications **31** (2020). https://doi.org/https://doi.org/10.1007/s00138-020-01088-z

39. Zare, A., Jiao, C., Glenn, T.: Discriminative multiple instance hyperspectral target characterization. IEEE transactions on pattern analysis and machine intelligence **40**(10), 2342–2354 (2017)

40. Zeng, G., Birchfield, S.T., Wells, C.E.: Detecting and measuring fine roots in minirhizotron images using matched filtering and local entropy thresholding. Machine Vision and Applications **17**(4), 265–278 (2006)

41. Zeng, G., Birchfield, S.T., Wells, C.E.: Rapid automated detection of roots in minirhizotron images. Machine Vision and Applications **21**(3), 309–317 (2010)

42. Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., Torralba, A.: Learning deep features for discriminative localization. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2921–2929 (2016)