# Accurate Detection in Volumetric Images using Elastic Registration based Validation

Dominic Mai[1,2], Jasmin Dürr[3], Klaus Palme[2,3], and Olaf Ronneberger[1,2]

[1]Computer Science Department, University of Freiburg
[2]BIOSS Centre of Biological Signalling Studies, University of Freiburg
[3]Institute for Biologie II, University of Freiburg
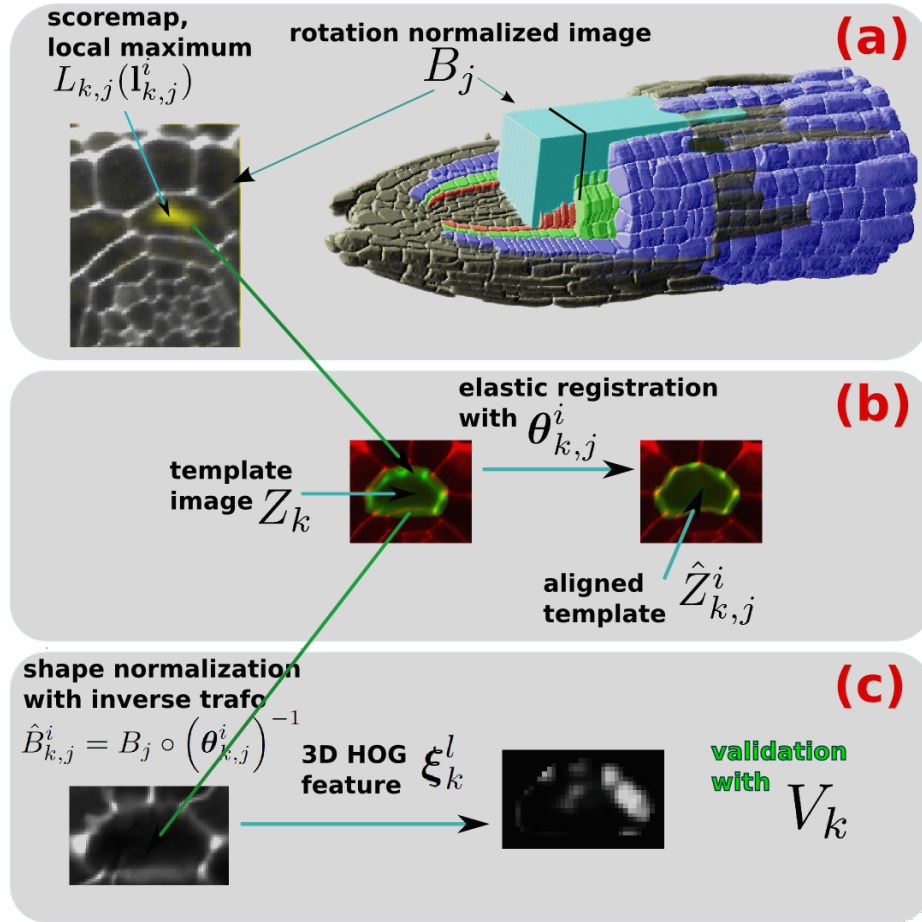`maid@informatik.uni-freiburg.de`

**Abstract.** In this paper, we propose a method for accurate detection and segmentation of cells in dense plant tissue of *Arabidopsis Thaliana*. We build upon a system that uses a top down approach to yield the cell segmentations: A discriminative detection is followed by an elastic alignment of a cell template. While this works well for cells with a distinct appearance, it fails once the detection step cannot produce reliable initializations for the alignment. We propose a validation method for the aligned cell templates and show that we can thereby increase the average precision substantially.

## 1 Introduction

Multi class segmentation is an important task in biomedical image analysis. It enables statistically meaningful analysis of signals by relating them to the underlying structures. There are basically two approaches to this problem: 1. In the bottom up approach, one generates a set of region hypotheses that are later classified and merged to obtain the class label, e.g. [9,3]. 2. In the top down approach one uses a detector to obtain coarse object localizations that are refined by a finer grained alignment of a model to the data, e.g. [10,2].

In [10] we presented a paper that deals with detection and alignment of plant cells in volumetric data in a top down approach. The goal of this paper was to detect single cells of a certain layer from an Arabidopsis root and to reconstruct this cell layer. *Arabidopsis thaliana* is a model organism widely used in plant biology [7,11]. We use a rigid cell detector based on 3D HOG features to coarsely localize the cells, similar to Dalal and Trigg's approach for 2D human detection [5]. Then we align a template image (*sharp mean image*) for the respective cell type to the data using elastic registration. Finally we reconstruct the root in a greedy fashion by assembling the aligned cell templates iteratively, beginning with the aligned detection whose associated detection filter received the highest score.

This approach works well for the cell layer 3 (Fig. 3), as the rigid detection filter produces reliable hypotheses for the alignment. Unfortunately it fails to produce satisfactory results for cell layer 4 as is illustrated in Fig. 2(a). The

**Fig. 1.** Overview of our processing pipeline. (a) We run a sliding window detector $D_k$ on the rotation normalized image $B_j$. On the left you see a slice of $B_j$ and a local maximum of the produced scoremap $L_{k,j}$. (b) The template image $Z_k$ is aligned to the data with the transformation $\boldsymbol{\theta}_{k,j}^i$. (c) We use the inverse transformation $(\boldsymbol{\theta}_{k,j}^i)^{-1}$ to shape normalize the root image at this location. Finally, we compute a 3D HOG feature and validate it with the proposed classifier $V_k$.

reason for this is the greedy reconstruction based on the scores obtained by the HOG based rigid detector. This coarse localization step is needed as it would be computationally impossible to perform alignments of all cell models in all image locations. The detection system is hence optimized to deliver a high recall. This, however, leads to many false positive detections, as the image data often allows for multiple different interpretations with similar scores. For example, it happens frequent that two adjacent small cells are interpreted as a bigger cell that encompasses both (or vice versa).

Therefore, the score from the rigid detector is a bad foundation to decide whether the suggested model describes the recorded data correctly. While the alignment step can correct for coarse localizations of the right cell type, it cannot correct the error if the wrong cell type has been chosen. This is especially bad with the greedy reconstruction approach: once a false alignment is accepted, it is likely to also prevent valid alignments in its direct neighborhood.

Bourdev et al. [1] use a linear support vector machine to rescore detections of persons based on mutually consistent poselet activations. Their framework, however, is more directed to deal with appearance changes due to the camera viewpoint and the articulated nature of their objects, opposed to the deformable nature of the plant cells that we consider.

*Contribution.* We propose an effective validation step that makes use of the finer grained localization of the elastic alignment. For every detector, we train a discriminative classifier that verifies whether the alignment of the sharp mean image is valid for a certain location. This validation step results in a much better greedy reconstruction (Fig. 2(b)).
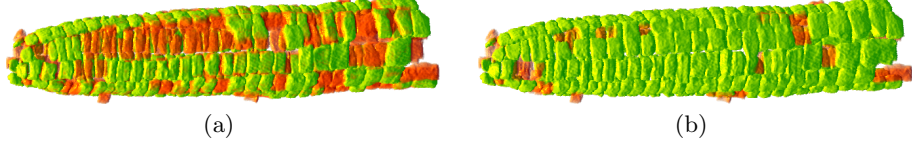
## 2   Detection and Alignment

The foundation of the approach presented here is our detection and alignment pipeline from [10]. We will outline the pipeline and introduce our formal notation along the way. For an overview, have a look at Fig. 1.

We define a 3D volumetric image $I$ of an Arabidopsis root as a function $I : \Omega \mapsto \mathbb{R}, \Omega \subset \mathbb{R}^3$. It comes with a set of ground truth cell segmentations $S_i : \Omega \mapsto \{0, 1\}$. Attached to the root is a root coordinate system (RCS) [12] consisting of the direction of the main axis of rotation of the root and an arbitrary but fixed "up" component perpendicular to this axis.

The root has a cylindrical structure with cells organized in concentric rings around its core (Fig. 3). The RCS is used to normalize for the orientation and the location of the cells with a rigid transformation $\mathbf{H}_i$. The normalized cells are clustered into $K$ clusters with a k-means clustering based on their cuboidal bounding volumes. For every cluster, a discriminative detection filter $D_k$ is trained. It is based on 3D HOG features that bin the image gradients into 20 orientation directions (vertices of a Dodecahedron). The soft binning and spatial pooling is realized with a convolution by a triangular filter with a radius $r_{\mathrm{HOG}}$. The HOG features are sampled on a regular grid at a distance of $s_{\mathrm{HOG}}$. The $D_k$ are realized as linear support vector machines, with the orientation normalized cell images of a cluster being the positive training examples and randomly sampled orientation normalized images from other parts of the root as negative examples. Along with the detector a *sharp mean image* $Z_k$ [13] is generated from the positive training examples. It represents the centroid of the cell cluster with respect to appearance and shape. The sharp mean image $Z_k$ comes with a segmentation mask $S_{Z_k}$.

At test time, all detection filters $D_k$ are tested in a sliding window fashion on overlapping rotation normalized cuboid shaped image regions

(a)                                              (b)

**Fig. 2.** Greedy reconstructions of layer 4 of the root r06. **(green)** cells are correctly aligned detections of cells with an IOU $\geq 0.5$, **(red)** cells are falsely aligned detections with an IOU $< 0.5$. (a) Reconstruction based on the scores from the rigid detector. Many locations of the root are occupied by false detections (average precision $= 0.64$). (b) The proposed validation step produces much better scores and thus leads to a better reconstruction (average precision $= 0.88$).

$B_j, j \in \{1, \ldots, N_B\}$ of the root, that are sampled in $10°$ steps. The rotation normalization is based on the RCS. The sliding window is realized as a convolution operation that is efficiently computed in the Fourier domain. This results in score maps $L_{k,j} : \mathbb{R}^3 \mapsto \mathbb{R}$. The detection locations $\mathbf{l}_{k,j}^i \in \mathbb{R}^3, i \in \{1, \ldots, N_{k,j}\}$ ($i$th detection for the detection filter $D_k$ on the rotation normalized image $B_j$) are the local maxima of these maps. All local maxima need to be $> 0$. Within the volume of the segmentation mask $S_{Z_k}$, all local maxima except the best scoring local maxima are suppressed.

The corresponding sharp mean image $Z_k : \mathbb{R}^3 \mapsto \mathbb{R}$ is put at the location $\mathbf{l}_{k,j}^i$ and is subsequently aligned to the rotation normalized image region $B_j$ with an elastic registration based on the combinatorial optimization from [8]. The elastic registration yields a transformation $\boldsymbol{\theta}_{k,j}^i : \mathbb{R}^3 \mapsto \mathbb{R}^3$ that is used for obtaining the aligned sharp mean image $\hat{Z}_{k,j}^i = Z_k \circ \boldsymbol{\theta}_{k,j}^i$ and the corresponding aligned segmentation mask $\hat{S}_{Z_k,j}^i = S_{Z_k} \circ \boldsymbol{\theta}_{k,j}^i$. Note that the $Z_k$ and $S_{Z_k}$ are only dependent on the $k$, but the aligned images $\hat{Z}_{k,j}^i$ and the aligned segmentation masks $\hat{S}_{Z_k,j}^i$ are also dependent on the respective rotation normalized image region $B_j$ and the implied transformation $\boldsymbol{\theta}_{k,j}^i$.

As last step, the aligned sharp mean images are transformed back into the coordinate system of the original root to obtain a reconstruction. This is done in an iterative greedy fashion, beginning with the aligned sharp mean image $\hat{Z}_{k,j}^i$ corresponding to the *highest scoring* detection location $\mathbf{l}_{k,j}^i$. The indices $\{k, j, i\}$ are formally given by

$$\{k, j, i\} = \arg \max_{\substack{k \in \{1, \ldots, K\} \\ j \in \{1, \ldots, N_B\} \\ i \in \{1, \ldots, N_{k,j}\}}} L_{k,j}(\mathbf{l}_{k,j}^i) \quad . \tag{1}$$

The aligned sharp mean image $\hat{Z}_{k,j}^i$ is transformed back to the original root image, then it is removed from the pool of available detection hypotheses and the detection hypotheses with the next best score is processed. Note that once a location in the original root image is occupied, it is not possible to put other aligned images at this location. This gives a crucial importance to the ordering

of the aligned candidates implied by Eqn. 1: Due to the continuous nature of the cells wrt. deformation, we usually have multiple competing aligned candidates per ground truth location. It is crucial to pick the well aligned candidates first during this greedy iterative reconstruction, as a badly aligned candidate can not be corrected and will probably also prevent subsequent well aligned candidates in its direct neighborhood.

If the scores delivered by the rigid detector fail to provide a good sorting of the aligned candidates prior to the greedy reconstruction, the results for layer 4 are not satisfactory (Fig. 2(a)). We propose a validation of the aligned sharp mean images that makes use of the finer grained information that is available due to the alignment. This results in much better reconstruction results as we will show in the experiments section (Fig. 2(b)).

## 3   Training of the Alignment Classifiers

In order to validate a candidate alignment of the sharp mean template we propose to use the metric induced by a discriminative classifier. To this end we will use a support vector machine, as it gives a normalized score around zero. Values $> 1$ indicate a confident decision for the positive class (well aligned candidate), values $< -1$ indicate a confident decision for the negative class (badly aligned candidate). As SVMs have good generalization properties, decision values in the interval $[-1, 1]$ mark a gradual change between the classes.

In our case the data used for training and testing is the 3D HOG representation of the root image within the support of the aligned cell template. The support vector machine, however, needs input data of a fixed size. This means that we cannot use the image data "below" the aligned template directly, as its volume is variable due to the elastic alignment. Therefore we will use the inverse transformation $\boldsymbol{\theta}^{-1}$ to warp the image data onto the cell template. This assures that all training and test data for a cluster $k$ will have the same number of features.

We need to mine positive $(+)$ and negative $(-)$ training examples from a training and validation root to train the validation classifier $V_k$. We compare the aligned segmentation masks $\hat{S}^i_{Z_k,j}$ with the ground truth segmentations. The *Intersection over Union* is the measure $M_{\mathrm{IOU}}$ used in the PASCAL VOC [6] challenge to assess the quality of a detection:

$$M_{\mathrm{IOU}}(S_1, S_2) = \frac{\int_{\Omega} S_1(\mathbf{x}) \cdot S_2(\mathbf{x}) \quad d\mathbf{x}}{\int_{\Omega} \max\left(S_1(\mathbf{x}), S_2(\mathbf{x})\right) \quad d\mathbf{x}} \quad . \tag{2}$$

This area based measure is well suited to evaluate the degree of alignment in a detection setting, especially when it is based on 3D segmentation masks. We assign the class $\{+, -\}$ based on the rule that is also used for the evaluation of the complete pipeline. An aligned candidate is accepted, iff the intersection over union of the corresponding aligned segmentation mask with ground truth

segmentation $S_l$ is greater than 0.5:

$$c(\hat{S}^i_{Z_k,j}) = \begin{cases} +, & M_{\text{IOU}}(\hat{S}^i_{Z_k,j}, S_l) \geq 0.5 \\ -, & M_{\text{IOU}}(\hat{S}^i_{Z_k,j}, S_l) < 0.5 \end{cases}. \tag{3}$$

We shape normalize the corresponding root image by transforming it with the inverse transformation:

$$\hat{B}^i_{k,j} = B_j \circ \left(\boldsymbol{\theta}^i_{k,j}\right)^{-1} \quad . \tag{4}$$

After the this transformation, we extract the 3D HOG feature that will be used in the training for the validation classifiers $V_k$.

$$\boldsymbol{\xi}^l_k = \mathbf{f}(\hat{B}^i_{k,j}), \quad \mathbf{f} : (\Omega \mapsto \mathbb{R}) \mapsto \mathbb{R}^{N^f_k} \tag{5}$$

The function $\mathbf{f}$ transforms the image $\hat{B}^i_{k,j}$ into a vectorial feature representation, i.e. the 3D HOG feature, and crops it along the support of the sharp mean image $Z_k$. For simplicity of notation we replace the indices $j, k$ with $l \in \{1, \ldots, N_k\}$, as its no longer important, from which $B_j$ the $\boldsymbol{\xi}^l_k$ originates. After the classification of the training examples into $(+)$ and $(-)$ with Eqn. 3 we end up with a set $\mathcal{S}^+_k = \{l^+_1, \ldots, l^+_{N^+}\}$ for positive examples and a set $\mathcal{S}^-_k = \{l^-_1, \ldots, l^-_{N^-}\}$ for the negative training examples for every cluster $k$.

As wish to investigate the effect of the model complexity of the classifier, train a *linear* support vector machine $V^{\text{lin}}_k$ a *RBF kernel* support vector machine $V^{\text{RBF}}_k$. We use 5-fold cross validation to estimate suitable parameters for the outlier penalty $c$ and the radius $\gamma$ of the radial basis function for $V^{\text{RBF}}_k$. The training is done with *libsvm* [4].
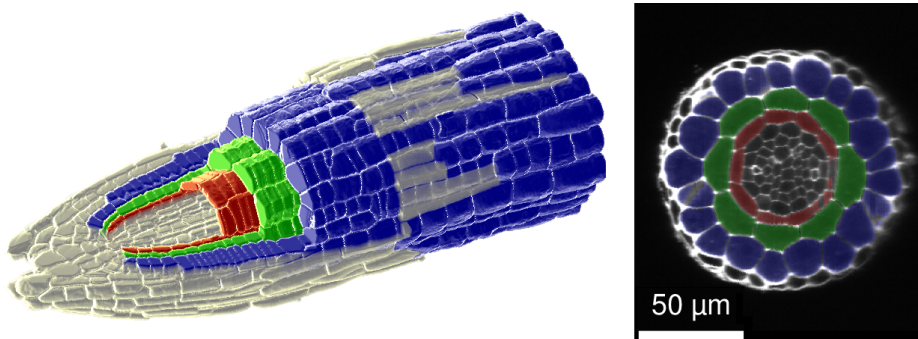
### 3.1   Validating Aligned Templates

The setting at test time is identical to the mining of the training examples for the $V_k$, except that we run the greedy iterative reconstruction of the root image after the detection and alignment phase. For each detection location $\mathbf{l}^i_{k,j}$, we perform the elastic registration of the corresponding sharp mean image $Z_k$ to the data and yield the transformation $\boldsymbol{\theta}^i_{k,j}$. Then we shape normalize the root image at this location by warping it with the inverse transformation $(\boldsymbol{\theta}^i_{k,j})^{-1}$ and compute the 3D HOG features $\boldsymbol{\xi}^l_k$ (Eqn. 5). Thus we end up with a set of aligned template image candidates and the corresponding 3D HOG feature representations of the locally shape normalized root image:

$$\left(\hat{Z}^l_k, \boldsymbol{\xi}^l_k\right) \text{with} \quad k \in \{1, \ldots, K\} \quad \text{and} \quad l \in \{1, \ldots, N_k\} \quad . \tag{6}$$

We perform the iterative greedy reconstruction, but replace the sorting induced by the rigid detector scores (Eqn. 1) with a sorting based on the proposed validation classifier. We begin with the best scoring candidate image

$$\hat{Z}^l_k \quad \text{with} \quad \{k, l\} = \arg \max_{\substack{k \in \{1, \ldots, K\} \\ l \in \{1, \ldots, N_k\}}} V^{\{\text{lin,RBF}\}}_k(\boldsymbol{\xi}^l_k) \quad . \tag{7}$$

**Fig. 3.** Volume rendering and slice of the raw data. The root has a cylindrical structure and is made up of concentric layers of different cell types. In this paper, we consider cells from **layer 2 (blue)**, **layer 3 (green)**, and **layer 4 (red)** for the detection task.

Note that we run the reconstruction either with the scores from the linear SVM $V_{\mathrm{lin}}$ or with the scores from the RBF SVM $V_{\mathrm{RBF}}$.

## 4    Experiments

In this section we show a quantitative and a qualitative evaluation of the effectiveness of the proposed validation approach. We had three Arabidopsis roots with ground truth segmentations available: r06, r14, and pi005. The generation of the ground truth is very time consuming, as each root contains $\sim$ 2500 cells. The ground truth segmentations are obtained by manually checking segmentations from a watershed algorithm on enhanced data [9]. For each root, we trained rigid detectors $D_k$ and validation classifiers $V_k$ for the cell layers 2, 3, and 4 (Fig. 3). We used a *round robin* scheme (Table 1), such that each root takes every role (training, validation, test) once. For the training of the rigid detectors we only used the *training* root. We always split the data into $k = 15$ clusters, as this value has proven to be good for layer 3 [10]. We do not perform a mining of hard negative examples for the detectors, as it did not improve the average precision of the results.

**Table 1.** Round robin scheme for training and testing.

| Training | Validation | Test |
|----------|-----------|-------|
| r06 | pi005 | r14 |
| r14 | r06 | pi005 |
| pi005 | r14 | r06 |

As the rigid detector returns virtually no false positives when tested on the training root, we mine the training examples for the validation classifier $V_k$ on the

*training* and the *validation* root. We train a linear support vector machine $V_k^{\text{lin}}$ and a RBF support vector machine $V_k^{\text{RBF}}$ using *libsvm* [4].

Our test setup is a detection setting. To assess the quality of the detections, we use the same method as in the PASCAL VOC challenge [6]. We accept a detection as valid, iff the intersection over union of the predicted segmentation mask with the ground truth segmentation is $\geq 0.5$. All subsequent detections of the same ground truth cell that are not suppressed during the reconstruction count as false positives. We investigate all combinations

$$\{\text{r06, r14, pi005}\} \times \{\text{layer 2, layer 3, layer 4}\} \times \{D_k, V_k^{\text{lin}}, V_k^{\text{RBF}}\}$$

and thus end up with 27 experiments. The sliding window detection is performed on the rotation normalized root images $B_j$ and takes $\sim 50s$ on a six core workstation. We compute the necessary convolutions in the Fourier domain, therefore the runtime is not dependent on the size of the detection filter $D_k$. The alignment of a cell template to the image data is computed with the combinatorial registration from our previous work [10], using a gradient orientation based data term. The computation of one alignment takes $\sim 1.5s$. The scoring of the aligned cell templates with the validation classifier takes $< 0.1s$. These steps are nearly perfectly parallelizable. When executed on a computing cluster with $5 \times 32$ cores the detection for a whole root takes $\sim 5$min, the alignment of the cell templates in average $\sim 20$min, depending on the number of cell hypotheses. The limiting factor with our setup was the hard disk IO.

The iterative greedy reconstruction is more difficult to parallelize and takes in average $\sim 30$min, also dependent on the number of aligned cell candidates. For some more statistics of the roots, have a look at Table 2.

**Table 2.** Statistics for the roots ("GT" = ground truth).

|  | **r06** | **r14** | **pi005** |
|---|---|---|---|
| size (voxels) | $1030 \times 433 \times 384$ | $944 \times 413 \times 360$ | $855 \times 458 \times 329$ |
| layer 2 #GT cells | 542 | 487 | 554 |
| layer 3 #GT cells | 216 | 188 | 208 |
| layer 4 #GT cells | 266 | 211 | 222 |
| $B_j$ arrangement | $3 \times 36$ | $3 \times 36$ | $2 \times 36$ |
| $B_j$ size (voxels) | $301 \times 101 \times 131$ | $301 \times 101 \times 131$ | $301 \times 101 \times 131$ |

Our findings are summarized in Fig. 4 as precision-recall graphs and in Table 3 as the mean average precision ($\varnothing$AP) per cell layer. The average precision is computed as the area under the precision-recall curve. Our original processing pipeline (cyan curve) works reasonably well for layer 2 and layer 3 with $\varnothing$AP $= 0.71$ and $\varnothing$AP $= 0.82$ respectively. It fails for layer 4 with $\varnothing$AP $= 0.52$. The reason for this can be found in the less distinctive cell shapes of this layer and in its location within the root. For the volumetric recording of layer 4, the light has to pass layers 1, 2, and 3 during the recording with a confocal microscope, which results in a more distorted signal.

When performing the greedy reconstruction based on the scores of the proposed validation approach, we yield substantially better results for the difficult layer 4. For an illustration see Fig. 2. For every other root and layer combination we also achieve better results through the validation scores. The linear SVM based scores (black curve) achieve the best reconstructions on layer 3. With the RBF SVM based scores (red curve), we achieve the best reconstructions for layer 2 and layer 4. The performance of the linear scoring and the RBF scoring is very similar, maybe with a slight edge for the RBF based rescoring. The training times of the SVMs are practically identical. When training directly on the kernel matrix with *libsvm* [4], the training takes under a minute including a cross validation based grid search for the SVM parameters $\gamma$ and $C$.

For the validation with the RBF SVM, one needs to compute between $50\times$ and $200\times$ longer compared to the validation with the linear SVM. However, the time needed to compute a validation is dominated by the disk IO, which leads to a similar overall computation time.

**Table 3.** Mean average precision ($\varnothing$AP) for scoring strategy and cell layer.
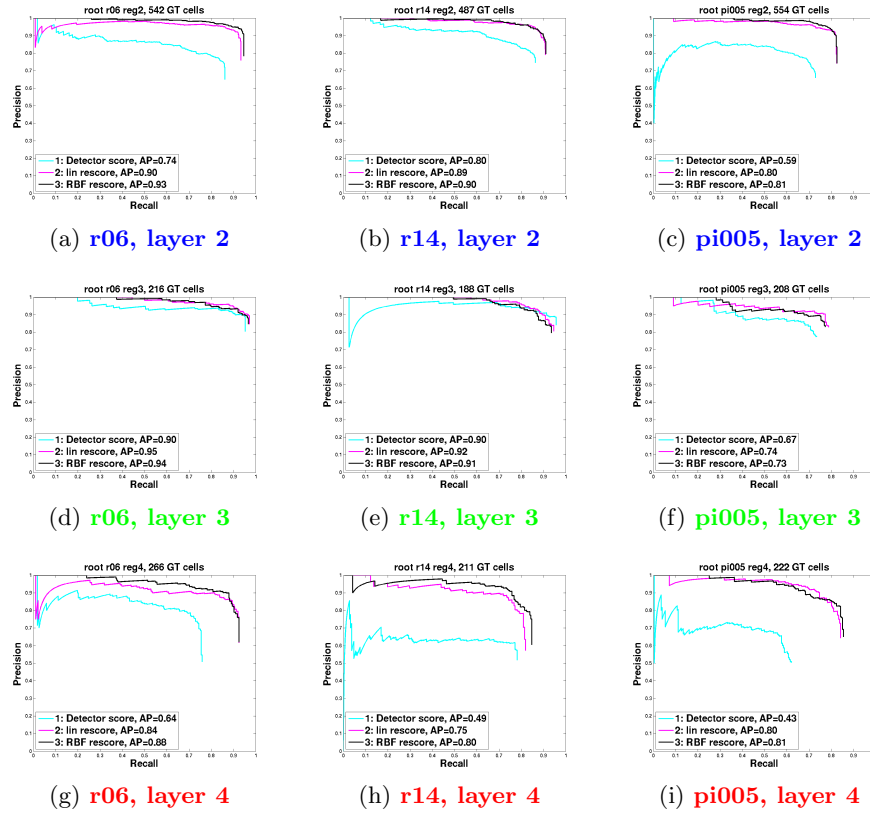
|  | layer 2 | layer 3 | layer 4 |
|---|---|---|---|
| raw detector score | 0.71 | 0.82 | 0.52 |
| linear SVM | 0.86 | **0.87** | 0.80 |
| RBF SVM | **0.88** | 0.86 | **0.83** |

## 5  Conclusions

In this paper we presented a validation strategy for detections in volumetric images that leverages the fine grained localization provided by the elastic alignment of a cell template image to the underlying image data. We use a metric based on trained discriminative classifiers to decide whether this alignment was successful or not. This validation step comes at practically no extra cost given the aligned detections. However, it achieves to boost the detection accuracy substantially, especially in regions of lower data quality, where the scores of the rigid detector are no longer reliable. We believe that this validation strategy should also work for other object classes in 2D and 3D images when the intra class appearance variation mainly stems from an elastic deformation of the objects.

## Acknowledgements

(a) **r06, layer 2**      (b) **r14, layer 2**      (c) **pi005, layer 2**

(d) **r06, layer 3**      (e) **r14, layer 3**      (f) **pi005, layer 3**

(g) **r06, layer 4**      (h) **r14, layer 4**      (i) **pi005, layer 4**

**Fig. 4.** Precision-Recall curves for cell layer reconstructions of the roots r06, r14, and pi005 organized in columns, e.g. root r06: (a), (d), (g) and cell layers 2, 3, and 4 organized in rows, e.g. layer 2: (a), (b), (c). **(cyan curve)** Reconstruction based on the detector scores. **(black curve)** Reconstruction based on the scores after validating the aligned cell templates with a linear SVM. **(magenta curve)** Reconstruction based on the scores after validating with an RBF SVM. All settings benefit from the better scores produced by the validation. The benefit is the biggest for layer 4, as the cells in this layer have the least distinctive cell shape and the data quality is worst for this layer. (d) is the configuration examined in our previous work [10].

## References

1. Bourdev, L., Maji, S., Brox, T., Malik, J.: Detecting people using mutually consistent poselet activations. In: ECCV (2010) 3
2. Brox, T., Bourdev, L., Maji, S., Malik, J.: Object segmentation by alignment of poselet activations to image contours. In: IEEE International Conference on Computer Vision and Pattern Recognition (CVPR) (2011) 1
3. Carreira, J., Sminchisescu, C.: CPMC: Automatic Object Segmentation Using Constrained Parametric Min-Cuts. IEEE Transactions on Pattern Analysis and Machine Intelligence (2012) 1

4. Chang, C.C., Lin, C.J.: LIBSVM: A library for support vector machines. ACM Transactions on Intelligent Systems and Technology 2, 27:1–27:27 (2011) 6, 8, 9
5. Dalal, N., Triggs, B.: Histograms of oriented gradients for human detection. In: Schmid, C., Soatto, S., Tomasi, C. (eds.) International Conference on Computer Vision & Pattern Recognition (2005) 1
6. Everingham, M., Van Gool, L., Williams, C.K.I., Winn, J., Zisserman, A.: The pascal visual object classes (voc) challenge. International Journal of Computer Vision 88(2), 303–338 (Jun 2010) 5, 8
7. Fernandez, R., Das, P., Mirabet, V., Moscardi, E., Traas, J., Verdeil, J., Malandain, G., Godin, C.: Imaging plant growth in 4d: robust tissue reconstruction and lineaging at cell resolution. Nature methods 7(7), 547–553 (2010) 1
8. Komodakis, N., Tziritas, G., Paragios, N.: Performance vs computational efficiency for optimizing single and dynamic mrfs: Setting the state of the art with primal-dual strategies. Computer Vision and Image Understanding 112(1), 14–29 (2008) 4
9. Liu, K., Schmidt, T., T.Blein, Dürr, J., Palme, K., Ronneberger, O.: Joint 3d cell segmentation and classification in the arabidopsis root using energy minimization and shape priors. In: IEEE International Symposium on Biomedical Imaging (ISBI) (2013) 1, 7
10. Mai, D., Fischer, P., Blein, T., Dürr, J., Palme, K., Brox, T., Ronneberger, O.: Discriminative detection and alignment in volumetric data. In: Weickert, J., Hein, M., Schiele, B. (eds.) GCPR. Lecture Notes in Computer Science, vol. 8142, pp. 205–214. Springer (2013) 1, 3, 7, 8, 10
11. Marcuzzo, M., Quelhas, P., Campilho, A., Maria Mendonça, A., Campilho, A.: Automated arabidopsis plant root cell segmentation based on svm classification and region merging. Comput. Biol. Med. 39(9), 785–793 (2009) 1
12. Schmidt, T., Pasternak, T., Liu, K., Blein, T., Aubry-Hivet, D., Dovzhenko, A., Dürr, J., Teale, W., Ditengou, F.A., Burkhardt, H., Ronneberger, O., Palme, K.: The irocs toolbox – 3d analysis of the plant root apical meristem at cellular resolution. The Plant Journal 77(5), 806–814 (Mar 2014), http://lmb.informatik.uni-freiburg.de//Publications/2014/SLBR14 3
13. Wu, G., Jia, H., Wang, Q., Shen, D.: Sharpmean: Groupwise registration guided by sharp mean image and tree-based registration. NeuroImage 56(4) (2011) 3