

# Joint Multi-Leaf Segmentation, Alignment, and Tracking for Fluorescence Plant Videos

Xi Yin, Xiaoming Liu, Jin Chen, David M. Kramer

**Abstract**—This paper proposes a novel framework for fluorescence plant video processing. The plant research community is interested in the leaf-level photosynthetic analysis within a plant. A prerequisite for such analysis is to segment all leaves, estimate their structures, and track them over time. We identify this as a joint multi-leaf **segmentation, alignment, and tracking** problem. First, leaf segmentation and alignment are applied on the last frame of a plant video to find a number of well-aligned leaf candidates. Second, leaf tracking is applied on the remaining frames with leaf candidate transformation from the previous frame. We form two optimization problems with shared terms in their objective functions for leaf alignment and tracking respectively. A quantitative evaluation framework is formulated to evaluate the performance of our algorithm with four metrics. Two models are learned to predict the alignment accuracy and detect tracking failure respectively in order to provide guidance for subsequent plant biology analysis. The limitation of our algorithm is also studied. Experimental results show the effectiveness, efficiency, and robustness of the proposed method.

**Index Terms**—plant phenotyping, Arabidopsis, leaf segmentation, alignment, tracking, multi-object, Chamfer matching



## 1 INTRODUCTION

PLANT phenotyping [1] refers to a set of methodologies and protocols used to measure plant growth [2], architecture [3], composition [4], and etc. In contrast to the manual observation-based methods, the automatic image-based approaches for plant phenotyping have gained more attention recently [5], [6]. As shown in Figure 1, plant researchers conduct large-scale experiments in a chamber with controlled temperature and lighting conditions. Ceiling-mounted fluorescence cameras capture images of a plant during its growth period [7]. The pixel intensities of the image indicate the *photosynthetic efficiency* (PE) of the plant. Given such a high-throughput imaging system, the massive data calls for advanced visual analysis in order to study a wide range of plant physiological problems [8], e.g., the heterogeneity of PE among the leaves. Therefore, the *leaf-level* visual analysis is fundamental to automatic plant phenotyping.

This paper focuses on the processing of the rosette plants where the leaves are at a similar height and form a circular arrangement. Our experiments are mainly conducted on *Arabidopsis thaliana*, which is the first plant to have its genome sequenced [9]. Due to its rapid life cycle, prolific seed production, and easiness to cultivate in the restricted space, *Arabidopsis* is the most popular and important model plant [10] in the plant research community. An automatic image analysis method for *Arabidopsis*, which is the main focus of this paper, is of essential importance for high-throughput plant phenotyping studies. Given a fluorescence plant video, our method performs multi-leaf Segmentation, Alignment, and Tracking (SAT) *jointly*. Specifically, leaf segmentation [11] segments each leaf

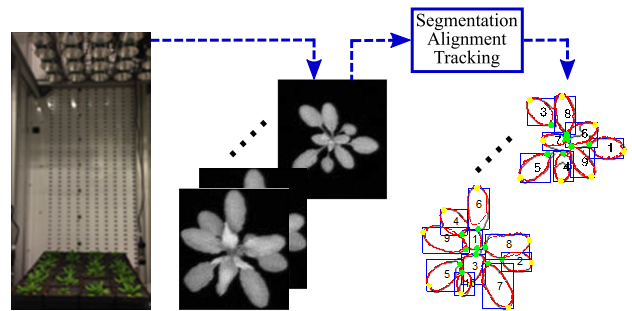


Fig. 1. Given a fluorescence plant video captured during its growth period, our algorithm performs multi-leaf SAT jointly, *i.e.*, estimating unique and consistent-over-time labels for all leaves and their individual leaf structure like leaf tips.

from the plant. Leaf alignment [12] estimates the leaf structure. Leaf tracking [13] associates the leaves over time. This multi-leaf analysis is a challenging problem due to several factors. First, the image resolution is low where the small leaves are even hard to be recognized by humans. Second, there are various degrees of overlap among leaves, which make it difficult to segment each leaf boundary. Third, leaves within a plant exhibit various shapes, sizes, and orientations, which also change over time. Therefore, an effective algorithm should be developed to handle all these challenges.

To the best of our knowledge, there is no previous work focusing on leaf SAT simultaneously from plant videos. To solve this new problem, we develop two optimization algorithms. Specifically, leaf segmentation and alignment are based on Chamfer Matching (CM) [14], which is a well-known algorithm to align one object in an image with a given template. However, classical CM does not work well for align-

ing multiple overlapping leaves. Motivated by crowd segmentation [15], where the number and locations of the pedestrians are estimated simultaneously, we propose a novel framework to jointly align multiple leaves in an image. First we generate a large set of leaf templates with various shapes, sizes, and orientations. Applying all templates to the edge map of a plant image leads to the same amount of transformed leaf templates. We adopt the local search method for optimization to select a subset of leaf candidates that can best explain the edge map of the test image.

While leaf segmentation and alignment work well for one image, applying it to every video frame independently does not enable tracking - associating aligned leaves over time. Therefore, we formulate leaf tracking on one frame as a problem of transforming multiple aligned leaf candidates from the previous frame. The tracking optimization initialized with results of the previous frame can converge very fast and thus results in enhanced leaf association and computational efficiency.

In order to estimate the alignment and tracking accuracy, two quality prediction models are learned respectively. We develop a quantitative analysis with four metrics to evaluate the multi-leaf SAT performance. Furthermore, the limitation of our algorithm is studied. In summary, we make four contributions:

- ◊ We identify a new computer vision problem of joint multi-leaf SAT from plant videos. We collect a dataset of *Arabidopsis* and make it publicly available.
- ◊ We propose two optimization algorithms to solve this multi-leaf SAT problem.
- ◊ We develop two quality prediction models to predict the alignment accuracy and tracking failure.
- ◊ We set up a quantitative evaluation framework to jointly evaluate the performance.

Compared to our earlier work [12], [16], we have made five main changes: 1) One term is modified in the tracking objective function. The proposed method is superior to [12], [16] on a larger dataset. 2) We develop two quality prediction models. 3) We enhance the performance evaluation procedure and add one metric to evaluate segmentation accuracy. 4) We study the limitation of our tracking algorithm and show its robustness to leaf template transformation. 5) We extend our method to apply on RGB images [5] and compare the segmentation results to [17].

## 2 PRIOR WORK

Plant image analysis has been studied in computer graphics [18]–[20] and computer vision [11], [21]. For example, a leaf shape and appearance model is proposed to render photo-realistic images of a plant [18]. A data-driven leaf synthesis approach is developed to produce realistic reconstructions of dense foliage [19]. These models may not be applied to fluorescence images due to the lack of leaf appearance information. There are prior computer vision work on tasks

such as leaf segmentation [11], [21], alignment [12], [22], tracking [13], [16], and identification [23]–[25]. However, most previous studies focus on only one or two of these tasks. In contrast, our method addresses three tasks of leaf SAT.

**Leaf Segmentation** can be classified into two categories: 1) segmentation of a detached leaf from natural [22], [26]–[29] or clean background [24]; 2) pixel-wise segmentation of each leaf from a plant [17], [25]. Methods in the first category are usually used as the first step for leaf classification or species identification. [24] uses pixel-based color classification for leaf segmentation from a white background. [25] proposes active contour deformation method for compound leaf segmentation and identification.

Our work belongs to the second category. It is very challenging due to leaf variation and overlapping. Tsaftaris et al. organized a collation study of leaf segmentation on rosette plants in 2015 [30], [31]. Our method is evaluated with other three methods. Two of them are based on superpixels and watershed transformation segmentation. [17] uses distance map-based leaf center detection and leaf split points detection for leaf segmentation.

**Leaf Alignment** aims to find the structure of a leaf, which is useful for leaf segmentation. [26] deforms a polygonal model to leaf shape fitting, where the base and tip points are used to define a leaf template. The same points are used in [22] to model leaf shapes and deform templates. Similarly, we use these two points on our leaf templates for alignment. Our novelty lies in solving leaf segmentation and alignment jointly by extending CM to align multiple potentially overlapping leaves in an image.

**Leaf Tracking** models leaf transformation over time. A probabilistic parametric active contour model is applied for leaf segmentation and tracking to automatically measure the temperature of leaves in [13]. However leaves on those images are well separated without any overlap and the active contours are initialized via the ground truth segments, which is hard to achieve in real-world applications. [32] segments all leaves in a video separately and employs a merging procedure to group the segments by exploiting the angle properties of the leaves. [33] proposes a graph-based tracking algorithm by linking leaf detections across neighboring frames. All of them treat tracking as a post processing *after* leaf segmentation on individual frame. In contrast, we employ a leaf template transformation to transfer the segmentation and alignment results between continuous frames.

## 3 OUR METHOD

Figure 2 shows our framework. Given a plant video, we first apply leaf segmentation and alignment on the last frame to generate a number of well-aligned leaf candidates. Leaf tracking is considered as an alignment problem with the leaf candidates initialized from

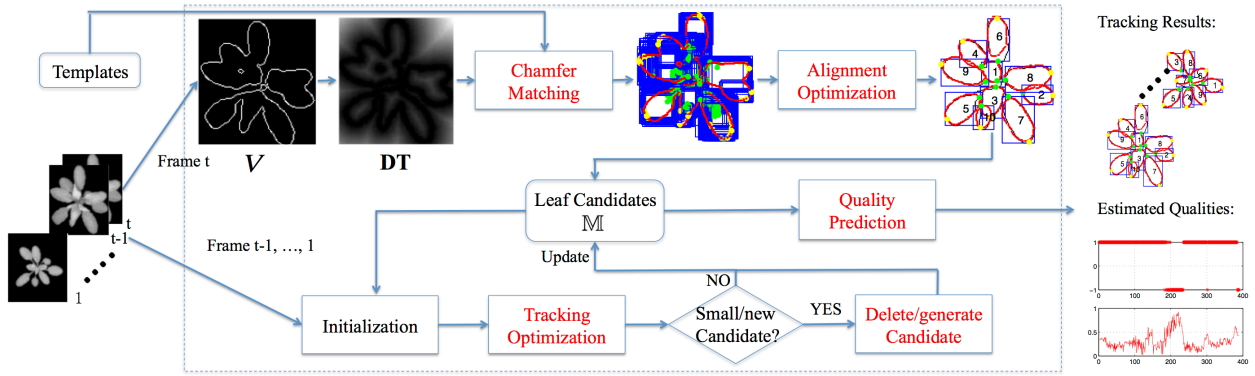


Fig. 2. Overview of the proposed joint multi-leaf SAT. Given a plant video with  $t$  frames, the proposed method outputs the SAT results on each frame, and two prediction curves on the quality of alignment and tracking for each leaf.

TABLE 1  
Notations.

Notation	Definition
$V, m$	the edge map and mask of a test image
$U, M$	the edge map and mask of a leaf template
$\tilde{U}, \tilde{M}$	the edge map and mask of a transformed leaf template
$DT$	the distance transform image of $V$
$H, S, R$	the numbers of template shapes, sizes, and orientations
$N$	the total number of leaf templates, $N = HSR$
$N_1$	the number of transformed leaf templates for optimization
$J, G$	the objective functions for alignment and tracking
$a$	the diagonal length of $V$
$(c^x, c^y)$	the center of a plant image
$(c_n^x, c_n^y)$	the center of the $n^{th}$ leaf candidate
$\mathbb{L}, \mathbb{L}_1$	the sets of $N$ and $N_1$ transformed templates
$A$	a $N_1 \times K$ matrix collecting all $\tilde{M}_n$ from $\mathbb{L}_1$
$K$	the number of pixels in the test image
$\mathbf{x}$	a $N_1$ -dim 0-1 indicator vector
$\mathbf{d}$	a $N_1$ -dim vector of CM distances in $\mathbb{L}_1$
$C$	a constant value used in $J_3$
$D$	the number of maximum iterations in tracking
$N^e, N^l$	the number of estimated and labeled leaves in a frame
$M$	the collection of $N^e$ selected leaf candidates
$P$	a set of transformation parameters $P = \{\mathbf{p}_n\}_{n=1}^{N^e}$ $\mathbf{p}_n = [\theta, r, t_x, t_y]^T$ is the parameter for $\tilde{U}_n$
$\hat{t}_{1,2}, t_{1,2}$	the estimated and labeled tips for one leaf
$\hat{T}, T$	the estimated and labeled tips for one frame
$\hat{\mathbb{T}}, \mathbb{T}$	the collections of estimated and labeled tips for all videos
$\hat{\mathbb{B}}, \mathbb{B}$	the collections of estimated and groundtruth segmentation masks for all videos
$N^b$	the total number of labeled leaves
$e_{1a}$	the tip-based error normalized by the leaf length
$ID$	a $N^e \times N^l$ matrix of leaf correspondence
$ER$	a $N^e \times N^l$ matrix of tip-based errors in one frame
$\mathbb{ER}$	the collection of all $ER$ for labeled frames
$f$	the number of leaf without correspondence
$e_1, e_2$	the tip-based errors used in Algorithm 2
$\tau$	a threshold for comparing with tip-based errors
$F, E, T$	the performance metrics
$Q_a, Q_t$	the quality to predict alignment and tracking
$\mathbf{x}_a, \mathbf{x}_t$	the features to learn quality prediction models
$\lambda_{1,2}, \mu_{1,2}$	the weights used in $J$ and $G$
$\alpha_1, \alpha_2$	the step sizes in the gradient descent of $J$ and $G$
$s$	the smallest leaf size we use

a previous frame. During tracking, a leaf candidate whose size is smaller than a threshold is deleted. A new candidate is detected and added for tracking when there is a certain region of the image mask that is not covered by the existing leaf candidates. Two prediction models are learned to investigate the alignment and tracking quality respectively. All notations are summarized in Table 1.

### 3.1 Multi-Leaf Segmentation and Alignment

Our segmentation and alignment algorithm consists of two steps. First, a pre-defined set of leaf templates is applied to the edge map of a test image to generate an over-complete set of transformed leaf templates. Second, we formulate an optimization process to select an optimal subset of leaf candidates.

#### 3.1.1 Candidate nomination via Chamfer matching

Chamfer Matching (CM) [14] is a well-known method used to find the best alignment between two edge maps. Let  $V = \{v_i\}$  and  $U = \{u_i\}$  be the edge maps of a test image and a template respectively. CM distance is computed as the average distance of each edge point in  $U$  with its nearest edge point in  $V$ :

$$d(U, V) = \frac{1}{|U|} \sum_{u_i \in U} \min_{v_j \in V} \|u_i - v_j\|_2, \quad (1)$$

where  $|U|$  is the number of edge points in  $U$ . CM distance can be computed efficiently via a pre-computed distance transform image  $DT(g) = \min_{v_j \in V} \|g - v_j\|_2$ , which calculates the distance of each coordinate  $g$  to its nearest edge point in  $V$ . During the CM process, an edge template  $U$  is superimposed on  $DT$  and the average value sampled by the template edge points  $u_i$  equals to the CM distance, i.e.,  $d(U, V) = \frac{1}{|U|} \sum_{u_i \in U} DT(u_i)$ .

Given a fluorescence plant image, it is first transformed to a binary image  $m$  by applying a threshold. The Sobel edge detector is applied to  $m$  to generate an edge map  $V$ . The goal of leaf alignment is to transform the 2D edge coordinates of a template  $U$  in the leaf template space to a new set of 2D coordinates  $\tilde{U}$  in the test image space so that the CM distance is small, i.e., the leaf template is well aligned with  $V$ .

**Image Warping:** In our framework, there are two types of transformations involved including forward and backward warping. We use affine transformation that consists of scaling, rotation, and translation.

As shown in Figure 3, let  $W : U \mapsto \tilde{U}$  be a *forward warping* function that transfers the 2D edge

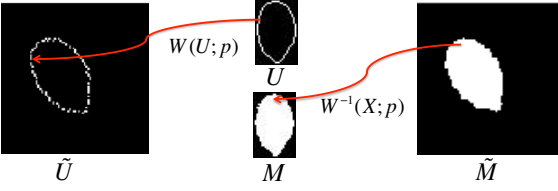


Fig. 3. Forward and backward warping.

points from the template space to the test image space, parameterized by  $\mathbf{p} = [\theta, r, t_x, t_y]^T$ :

$$\tilde{U} = \mathbf{W}(\mathbf{U}; \mathbf{p}) = r \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} (\mathbf{U} - \bar{U}) + \begin{bmatrix} t_x \\ t_y \end{bmatrix} + \bar{U}, \quad (2)$$

where  $\theta$  is the in-plane rotation angle,  $r$  is the scaling factor,  $t_x$  and  $t_y$  are the translations along  $x$  and  $y$  axis respectively.  $\bar{U}$  is the center of the leaf, *i.e.*, the average of all coordinates of  $U$ , which is used to model the leaf scaling and rotation w.r.t. the leaf center.

Let  $\mathbf{W}^{-1} : \tilde{U} \mapsto U$  be the *backward warping* from the image space to the template space. We denote  $\mathbf{X}$  as a  $K \times 2$  matrix including all coordinates in the test image space. Thus,  $\mathbf{W}^{-1}(\mathbf{X}; \mathbf{p})$  are the corresponding coordinates of  $\mathbf{X}$  in the template space. The purpose for this backward warping is to generate a  $K$ -dim vector  $\tilde{M} = M(\mathbf{W}^{-1}(\mathbf{X}; \mathbf{p}))$ , which is the warped version of the original template mask  $M$ .

**Leaf Templates:** Since there is a large variation in leaf shapes, it is infeasible to match leaf with one template. We manually select  $H$  basic templates (the 1<sup>st</sup> row in Figure 4) with representative leaf shapes from the plant videos and compute their individual edge map  $U$  and mask  $M$ . We synthesize an over-complete set of transformed templates by selecting a discrete set of  $\theta$  and  $r$ , which are expected to cover all potential leaf configurations in  $\mathbf{V}$ . This leads to an array of  $N = HSR$  leaf templates where  $S$  and  $R$  are the numbers of leaf scales and orientations respectively (Figure 4). The yellow and green points in Figure 4 are the two labeled leaf tips  $\mathbf{t}$ , which are used to find the corresponding leaf tips  $\hat{\mathbf{t}}$  in  $\mathbf{V}$  via Equation 1.

For each template  $U_n$ , it scans through all possible locations on  $\mathbf{V}$  and the location with the minimum CM distance is selected, which provides  $t_x$  and  $t_y$  optimal to  $U_n$ . Therefore, with the manually selected  $\theta, r$ , and exhaustively chosen  $t_x$  and  $t_y$ ,  $N$  transformed templates are generated from  $H$  basic templates. For each transformed template, we record the 2D edge coordinates of its basic template, warped template mask, transformation parameters, CM distance and the estimated leaf tips as  $\mathbb{L} = \{U_n, \tilde{M}_n, \mathbf{p}_n, d_n, \hat{\mathbf{t}}_n\}_{n=1}^N$ . Note that  $\mathbb{L}$  is an over-completed set of transformed leaf templates including the true leaf candidates as its subset. Hence, the critical question is how to select such a subset of candidates from  $\mathbb{L}$ .

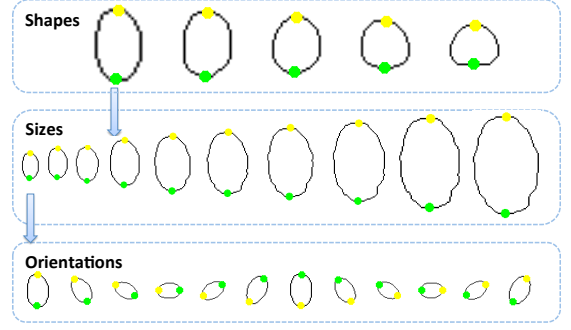


Fig. 4. Leaf template scaling and rotation from basic template shapes. The tip labels are shown in yellow and green.

### 3.1.2 Objective function

The goal of leaf segmentation and alignment is to segment each leaf and estimate the structure precisely. If the leaf candidates are well selected, there should be no redundant or missing leaves. Each leaf candidate should be well aligned with the edge map of the test image. This rationality leads to a three-term objective function, which seeks the minimal number of leaf candidates ( $J_1$ ) with small CM distances ( $J_2$ ) to best cover the test image mask ( $J_3$ ).

Each image contains around 10 leaves while the number of potential candidates in  $\mathbb{L}$  is 2,880 in our case. The selection space needs to be narrowed down substantially. To do this, we compute the CM distance and the overlap ratio of each template to the test image mask. We remove leaf templates whose CM distance is larger than the average of all templates or whose overlap ratio is smaller than 90%. Finally, we generate a new set  $\mathbb{L}_1$  with  $N_1$  (a few hundreds) templates. RANSAC [34] is not applicable here for two reasons. First, it is difficult to define a model or evaluation criterion for a random subset of leaf templates. Second, we have more outliers than inliers, which makes it hard to select the correct set in consensus.

The objective function is defined on a  $N_1$ -dim indicator vector  $\mathbf{x}$ , where  $x_n = 1$  means that the  $n^{\text{th}}$  transformed template is selected and  $x_n = 0$  otherwise. Hence  $\mathbf{x}$  uniquely specifies a combination of transformed templates from  $\mathbb{L}_1$ . The first term is the number of the selected leaf candidates  $J_1 = \|\mathbf{x}\|_1$ .

We concatenate  $d_n$  from  $\mathbb{L}_1$  to form a  $N_1$ -dim vector  $\mathbf{d}$ . The second term, *i.e.*, the average CM distance of the selected leaf candidates, is formulated as:

$$J_2 = \frac{\mathbf{d}^T \mathbf{x}}{\|\mathbf{x}\|_1}. \quad (3)$$

The third term is the comparison between the synthesized mask and the test image mask. As shown in Figure 5, we convert the binary mask to a  $K$ -dim row vector  $\mathbf{m}$  by raster scan. Similarly, each warped template mask  $\tilde{M}_n$  is also a  $K$ -dim vector. The collection of  $\tilde{M}_n$  from all transformed templates is denoted as a  $N_1 \times K$  matrix  $\mathbf{A}$ . Note that  $\mathbf{x}^T \mathbf{A}$  is indicative of the synthesized mask except that the

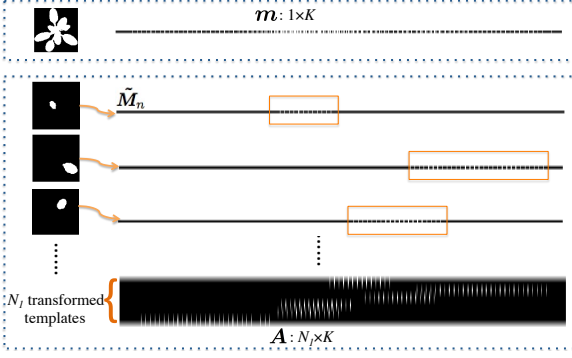


Fig. 5. The process of generating  $m$  and  $A$ .

pixel values of the overlapping leaves are larger than 1. We employ the  $\arctan(\cdot)$  function, similar to [35], [36], to convert all elements in  $x^T A$  to be in the range of  $[0, 1]$ ,

$$f(x) = \frac{1}{\pi} \arctan(C(x^T A - \frac{1}{2})) + \frac{1}{2}, \quad (4)$$

where  $C$  is a constant controlling how close  $\arctan(\cdot)$  approximates the step function. Note that the actual step function cannot be used here since it is not differentiable and thus is difficult to optimize. The constant  $\frac{1}{2}$  within the parentheses is a flip point separating where the value of  $x^T A$  will be pushed toward either 0 or 1. Therefore, the third term becomes:

$$J_3 = \frac{1}{K} \|f(x) - m\|_2^2. \quad (5)$$

Finally, our objective function is:

$$J(x) = J_1 + \lambda_1 J_2 + \lambda_2 J_3, \quad (6)$$

where  $\lambda_1$  and  $\lambda_2$  are the weights. These three terms jointly provide guidance on what constitutes an optimal combination of leaf candidates.

### 3.1.3 Local search method for optimization

Equation 6 is a *pseudo-Boolean function*. The basic algorithm [37] is not applicable because our objective cannot be written in the required polynomial form. We adopt the widely used local search method to optimize Equation 6. The local search algorithm [38] for pseudo-Boolean function iteratively searches a small neighborhood of  $x$  and updates  $x$  to its neighborhood that leads to a smaller function value.

First, all elements in  $x$  are initialized as 1, *i.e.*, all transformed templates are selected. We fix one element in  $x$  at each iteration by searching the neighborhood of  $x$  with the  $n^{\text{th}}$  element being 0 or 1, denoted as  $x_{x_n=0}$  and  $x_{x_n=1}$ . According to the proposition 6 in [38], a positive gradient indicates 0 in the corresponding element of the local optimal solution. Therefore, each iteration, we select the element with the maximum gradient to remove redundant leaf

templates. The gradient of the objective w.r.t.  $x$  is:

$$\begin{aligned} \frac{dJ}{dx} &= \text{sign}(x) + \lambda_1 \left( \frac{d}{\|x\|_1} - \frac{d^T x}{\|x\|_1^2} \text{sign}(x) \right) \\ &- \frac{2\lambda_2 C}{\pi K} A \left[ \left( f(x) - m \right) \odot \left( 1 + \left( C(x^T A - \frac{1}{2}) \right)^2 \right) \right]^T, \end{aligned} \quad (7)$$

where  $\text{sign}(x)$  is a function returning the sign of each element, and  $\odot$  is the element-wise division of vectors. In each iteration,  $x$  is updated by  $x = x - \alpha \frac{dJ}{dx}$ . The element  $x_n$  with the largest gradient is chosen and fixed to be 0 or 1 based on the smaller value of  $J(x_{x_n=0})$  and  $J(x_{x_n=1})$ . Once this element is fixed, its value remains unchanged in the future iterations. The total number of iterations is the number of transformed leaf templates  $N_1$ . Finally, those elements in  $x$  equal to 1 provide the combination of leaf candidates.

This joint leaf segmentation and alignment is applied on the last frame of a plant video to generate  $N^e$  leaf candidates that are used for tracking in the remaining video frames. We denote the set of leaf candidates selected from  $\mathbb{L}_1$  as  $\mathbb{M} = \{U_n, \tilde{M}_n, p_n, d_n, \hat{t}_n\}_{n=1}^{N^e}$ , which means the basic leaf template  $U_n$  is transformed by  $p_n$  to result in a leaf candidate that is well-aligned with the edge map.

## 3.2 Multi-Leaf Tracking Algorithm

Leaf tracking aims to assign the same leaf ID to the same leaf through an entire video. In order to track all leaves over time, one way is to apply leaf segmentation and alignment framework on every frame of the video and then build leaf correspondence between consecutive frames. However, the leaf tracking consistency is an issue due to the potentially different leaf segmentation results on different frames. Therefore, we form an optimization problem for leaf tracking based on template transformation.

### 3.2.1 Objective function

Similar to Equation 6, we formulate a three-term objective function parameterized by a set of transformation parameters  $P = \{p_n\}_{n=1}^{N^e}$ , where  $p_n$  is the transformation parameters for leaf candidate  $U_n$ .

First,  $P$  is updated so that the transformed leaf candidates are well aligned with the edge map  $V$ . The first term is computed as the average CM distance of the transformed leaf candidates:

$$G_1 = \frac{1}{N^e} \sum_{n=1}^{N^e} d(W(U_n; p_n), V). \quad (8)$$

The second term is to encourage the synthesized mask from all transformed candidates to be similar to the test mask  $m$ . The synthesized mask of one transformed leaf candidate is  $M_n(W^{-1}(X; p_n))$ , we formulate the second term as:

$$G_2 = \frac{1}{K} \left\| \sum_{n=1}^{N^e} M_n(W^{-1}(X; p_n)) - m \right\|_2^2. \quad (9)$$

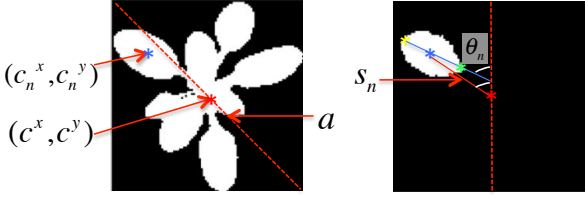


Fig. 6. The angle difference - the long axis of leaves should point to the plant center.

One property of rosette plants such as *Arabidopsis* is that the long axes of most leaves point toward the center of the plant. To take advantage of this domain-specific knowledge, the third term encourages the rotation angle  $\theta$  to be similar to the direction of the leaf center to the plant center. Figure 6 shows the geometric relation of the angle difference, which can be computed as  $\frac{c_n^x - c_x}{s_n} - \sin \theta_n$ , where  $(c^x, c^y)$  and  $(c_n^x, c_n^y)$  are the geometric centers of a plant and a leaf, i.e., the average coordinates of all points in  $\mathbf{m}$  and  $\bar{\mathbf{U}}_n$  respectively,  $s_n = \sqrt{(c_n^x - c_x)^2 + (c_n^y - c_y)^2}$  is the distance between the leaf center and the plant center, and  $\theta_n$  is the rotation angle. Furthermore, since this property is more dominant for leaves far away from the plant center, we weight the above angle difference by  $s_n$  and normalize it by the image size. The third term is the average weighted angle difference:

$$G_3 = \frac{1}{N^e a^2} \sum_{n=1}^{N^e} \|(c_n^x - c_x) - s_n \sin \theta_n\|_2^2. \quad (10)$$

Finally, the objective function is formulated as:

$$G(\mathbf{P}) = G_1 + \mu_1 G_2 + \mu_2 G_3, \quad (11)$$

where  $\mu_1$  and  $\mu_2$  are the weights.

Note the differences in two objective functions  $J$  and  $G$ . Since the number of leaves is fixed for tracking,  $J_1$  is not needed in the formulation of  $G$ . The number of leaves is relatively small during tracking. Therefore,  $\arctan(\cdot)$  is not needed since the synthesized mask is already comparable to the test image mask.

### 3.2.2 Gradient descent optimization

Given the objective function in Equation 11, our goal is to minimize it by estimating  $\mathbf{P}$ , i.e.,  $\mathbf{P} = \arg\min_{\mathbf{P}} G(\mathbf{P})$ . Since  $G(\mathbf{P})$  involves texture warping, it is a nonlinear optimization problem without a closed-form solution. We use gradient descent to solve this problem. The derivation of  $G_1$  w.r.t.  $\mathbf{P}$  can be written as:

$$\frac{dG_1}{d\mathbf{p}_n} = \frac{1}{N^e |\mathbf{U}_n|} (\nabla \mathbf{DT}^x \frac{\partial \mathbf{W}^x}{\partial \mathbf{p}_n} + \nabla \mathbf{DT}^y \frac{\partial \mathbf{W}^y}{\partial \mathbf{p}_n}), \quad (12)$$

where  $\nabla \mathbf{DT}^x$  and  $\nabla \mathbf{DT}^y$  are the gradient images of  $\mathbf{DT}$  at  $x$  and  $y$  axis respectively. These two gradient images only need to be computed once for each frame.  $\frac{\partial \mathbf{W}^x}{\partial \mathbf{p}_n}$  and  $\frac{\partial \mathbf{W}^y}{\partial \mathbf{p}_n}$  can be easily computed from Equation 2 w.r.t.  $\theta$ ,  $r$ ,  $t_x$  and  $t_y$  separately.

Similarly, the derivation of  $G_2$  w.r.t.  $\mathbf{P}$  is:

$$\frac{dG_2}{d\mathbf{p}_n} = \frac{2}{K} \left[ \sum_{n=1}^{N^e} M_n(\mathbf{W}^{-1}(\mathbf{X}; \mathbf{p}_n)) - \mathbf{m} \right] \cdot (\nabla M_n^x \frac{\partial \mathbf{W}_x^{-1}}{\partial \mathbf{p}_n} + \nabla M_n^y \frac{\partial \mathbf{W}_y^{-1}}{\partial \mathbf{p}_n}), \quad (13)$$

where  $\nabla M_n^x$  and  $\nabla M_n^y$  are the gradient images of the template mask  $M_n$  at  $x$  and  $y$  axis respectively.  $\frac{\partial \mathbf{W}_x^{-1}}{\partial \mathbf{p}_n}$  and  $\frac{\partial \mathbf{W}_y^{-1}}{\partial \mathbf{p}_n}$  can be computed based on the inverse function of Equation 2.

The derivation of  $G_3$  w.r.t.  $\theta$  is more complex than to the other three transformation parameters. For clarity, we only present the derivative over  $\theta$ :

$$\frac{dG_3}{d\theta_n} = \frac{2}{N^e a^2} [(c_n^x - c_x) - s_n \sin \theta_n] \cdot \left[ \frac{1}{|\mathbf{U}_n|} \frac{\partial \mathbf{W}_x}{\partial \theta_n} - s_n \cos \theta_n - \frac{2 \sin \theta_n}{s_n |\mathbf{U}_n|} [(c_n^x - c_x) \frac{\partial \mathbf{W}_x}{\partial \theta_n} + (c_n^y - c_y) \frac{\partial \mathbf{W}_y}{\partial \theta_n}] \right]. \quad (14)$$

During optimization,  $\mathbf{P}^0$  is initialized as the transformation parameters of the leaf candidates from the previous frame and updated as  $\mathbf{p}_n^t = \mathbf{p}_n^{t-1} - \alpha_2 \frac{dG}{d\mathbf{p}_n}$  for each leaf at iteration  $t$ . Note that this is a multi-leaf *joint* optimization problem because the computation of  $\frac{dG_2}{d\mathbf{p}_n}$  involves all  $N^e$  leaf candidates. The optimization stops when  $G$  does not decrease or it reaches the maximum iteration  $D$ .

### 3.2.3 Leaf candidates update

Given a multi-day plant video, we apply leaf segmentation and alignment algorithm on the last frame to generate  $\mathbb{M}$  and employ the leaf tracking toward the first frame. Due to plant growth and leaf occlusion, the number of leaves may vary throughout the video. If the area of any leaf candidate at one frame is less than a threshold  $s$  (defined as the number of pixels), we remove it from the leaf candidates.

On the other hand, a new leaf candidate can be detected and added to  $\mathbb{M}$ . To do this, we compute the synthesized mask of all leaf candidates and subtract it from the test image mask  $\mathbf{m}$  to generate a residue image for each frame. Connected component analysis is applied to find components that are larger than  $s$ . We then apply a subset of  $N$  leaf templates to find a leaf candidate based on the edge map of the residue image. The new candidate is assigned with an existing leaf ID if its overlap to a previous disappeared leaf is larger than a threshold. Otherwise it will be assigned with a new leaf ID. The new candidate is added into  $\mathbb{M}$  and tracked in the remaining frames.

## 3.3 Quality Prediction

While many algorithms strive for perfect results, it is inevitable that unsatisfactory or failed results are obtained on the challenging samples. It is critical for an algorithm to be aware of this situation so that

future analysis does not rely on poor results. One approach to achieve this goal is to perform the *quality prediction* for the task, similar to quality estimation for fingerprint [39] and face [40]. The key tasks in our work include leaf alignment, estimating the two tips of a leaf, and leaf tracking, keeping leaf consistency over time. Therefore, we learn two quality prediction models to predict the alignment accuracy and detect the tracking failure respectively. The prediction can be used to select a subset of leaves with high *quality* for subsequent plant biology analysis [41].

### 3.3.1 Alignment quality

Suppose  $Q_a$  is the alignment accuracy of a leaf, which indicates how well the two tips are aligned. We envision what factors may influence the estimation of the two tips. First, the CM distance indicates how well the template and the test image are aligned. Second, a well-aligned leaf candidate should have large overlap with the test image mask and small overlap with the neighboring leaves. Third, the leaf area, angle, and distance to the plant center may influence the alignment result. Therefore, we extract a 6-dim feature vector  $\mathbf{x}_a$  including: the CM distance  $d(\mathbf{W}(\mathbf{U}_n; \mathbf{p}_n), \mathbf{V})$ , the overlap ratio with the test image mask  $\frac{1}{|\mathbf{m}|_1} \tilde{\mathbf{M}}_n \odot \mathbf{m}$ , the overlap ratio with the other leaves  $\frac{\tilde{\mathbf{M}}_n \odot (\mathbf{m} - \tilde{\mathbf{M}}_n)}{|\tilde{\mathbf{M}}_n|_1}$ , the area normalized by test image mask  $\frac{1}{|\mathbf{m}|_1} |\tilde{\mathbf{M}}_n|_1$ , the angle difference  $|\theta_n - \sin^{-1} \frac{c_n^x - c_n^y}{s_n}|$  and the distance to the plant center  $s_n$ . A linear regression model is learned by optimizing the following objective on  $N^a$  training leaves with ground truth  $Q_a$ , which is proportional to the alignment error (details in Sec. 5.3.3).

$$\omega = \arg \min_{\omega} \sum_{n=1}^{N^a} \|Q_a^n - \omega \mathbf{x}_a^n\|, \quad (15)$$

where  $\omega$  is a 6-dim weighting vector to predict the alignment accuracy of each leaf.

### 3.3.2 Tracking quality

Due to the limitation of our algorithm, it is possible that one leaf might diverge to the location of the adjacent leaves and results in tracking inconsistency. We name it as a *tracking failure*. One example is shown in Figure 7, where labeled leaf 1 has been assigned two different IDs (3 and 4) during tracking. The change happens from frame 3 to frame 2. The goal of tracking quality prediction is to detect the *moment* when tracking starts to fail. We denote tracking quality as  $Q_t$ , where  $Q_t = -1$  means a tracking failure of one leaf and  $Q_t = 1$  means tracking success.

Similar to Section 3.3.1, we first extract a 6-dim feature vector  $\mathbf{x}_a$  for one leaf. However  $\mathbf{x}_a$  alone can not predict the tracking failure because it does not include temporal information. So we compare the features  $\mathbf{x}_a$  of one frame with that of a reference frame  $\hat{\mathbf{x}}_a$ , which is 20 frames before  $\mathbf{x}_a$ . Since a tracking

failure may result in abnormal changes in leaf area, angle, and distance to the center, we compute the leaf angle difference, leaf center distance, leaf overlap ratio between the current and the reference frame. Finally, we form a 15-dim feature vector denoted as  $\mathbf{x}_t$ , including  $\mathbf{x}_a$ ,  $\mathbf{x}_a - \hat{\mathbf{x}}_a$ , the leaf angle difference  $\theta_n - \hat{\theta}_n$ , the leaf center distance  $\sqrt{(c_n^x - \hat{c}_n^x)^2 + (c_n^y - \hat{c}_n^y)^2}$ , and the leaf overlap ratio  $\frac{\mathbf{M}_n(\mathbf{W}^{-1}(\mathbf{X}; \mathbf{p}_n)) \odot \tilde{\mathbf{M}}_n(\mathbf{W}^{-1}(\mathbf{X}; \hat{\mathbf{p}}_n))}{\mathbf{M}_n(\mathbf{W}^{-1}(\mathbf{X}; \mathbf{p}_n))}$ . Given a training set  $\Omega = \{(\mathbf{x}_t^n, Q_t^n)\}$  with  $Q_t^n \in \{-1, 1\}$ , a SVM classifier is learned as the tracking quality model.

## 4 PERFORMANCE EVALUATION

Leaf segmentation is to segment each leaf from the image. Leaf alignment is to correctly estimate two tips of each leaf. Leaf tracking is to keep the leaf ID consistent over the video. In order to quantitatively evaluate the performance of joint multi-leaf SAT, we need to provide the ground truth of the pixel-level leaf segments in each frame, the two tips of each leaf, and the leaf IDs for all leaves in the video.

As shown in Figure 7, we label the two tips of each leaf and manually assign their IDs in several frames of one video. We record the label results in one frame as a  $N^l \times 4$  matrix  $\mathbf{T}$ , where  $N^l$  is the number of labeled leaves and  $\mathbf{T}(n, :) = [t_1^x, t_1^y, t_2^x, t_2^y]$  records tip coordinates of  $n^{\text{th}}$  leaf in this frame. The collection of all labeled frames in all videos is denoted as  $\mathbb{T}$ , where  $\mathbf{T} = \mathbb{T}\{i, j\} (i = 1, 2, \dots, m, j = 1, 2, \dots, n)$ ,  $m$  is the number of labeled videos and  $n$  is the number of labeled frames in each video. The total number of labeled leaves in  $\mathbb{T}$  is  $N^b$ .

During template transformation, the corresponding points of the transformed template tips in  $\mathbf{V}$  become the estimated leaf tips  $[\hat{t}_1^x, \hat{t}_1^y, \hat{t}_2^x, \hat{t}_2^y]$ . The leaf ID is assigned in the last frame starting from 1 to the total number of selected leaves and kept the same during tracking. Similar to the data structure of  $\mathbb{T}$ , the tracking results of all videos over the labeled frames is written as  $\hat{\mathbb{T}}$ . Given  $\hat{\mathbb{T}}$  and  $\mathbb{T}$ , Algorithm 2 provides our detailed performance evaluation, which is also illustrated by a synthetic example in Figure 7.

There are two concepts involved: frame-to-frame and video-to-video correspondence. For each estimated leaf, we need to find one corresponding leaf in the labeled frame. Frame-to-frame correspondence aims to assign a unique leaf ID to each leaf in the frame so that the IDs are consistent with our labeled IDs. As mentioned before, the frame-to-frame correspondence may not be consistent in the whole video due to the tracking failures or more than one leaf IDs can be assigned to the same leaf in the video. Video-to-video correspondence aims to assign consistent and unique leaf ID to the same leaf in the entire video.

We start by building frame-to-frame leaf correspondence, as in Algorithm 1 and the red dotted box in Figure 7. To build the leaf correspondence of  $N^e$

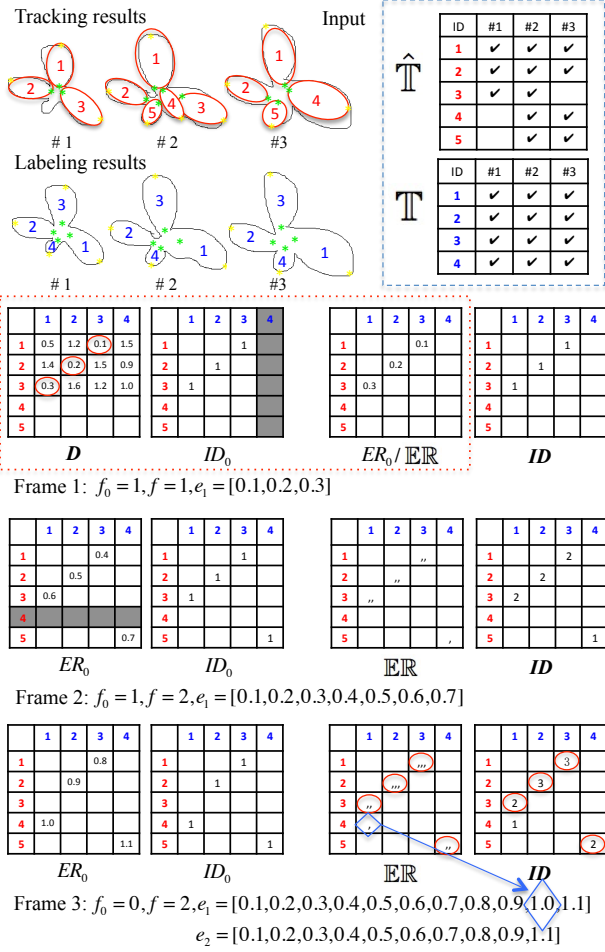


Fig. 7. A toy example of executing Step 1 of Algorithm 2 on one video with 3 frames. In frame 1, we illustrate the process of Algorithm 1. Each table shows the corresponding computation of each leaf from tracking results (each row) and label results (each column).

**Algorithm 1:** Build leaf correspondence  $[f, ER, ID] = \text{leafMatch}(\hat{T}, T)$ .

**Input:** Estimated leaf tips matrix  $\hat{T} (N^e \times 4)$  and labeled leaf tips matrix  $T (N^l \times 4)$ .

**Output:**  $f = |N^l - N^e|$ ,  $ER$ , and  $ID$ .

Initialize  $D = ER = ID = \mathbf{0}_{N^e \times N^l}$ .

**for**  $i = 1, \dots, N^e$  **do**

**for**  $j = 1, \dots, N^l$  **do**  
          $D(i, j) = e_{la}(\hat{t}_i, \hat{t}_j)$ ;

**for**  $k = 1, \dots, \min(N^e, N^l)$  **do**

$[e_{min}, i, j] = \min(D)$ ;  
      $D(i, :) = \text{Inf}$ ;  $D(:, j) = \text{Inf}$ ;  
      $ID(i, j) = 1$ ;  $ER(i, j) = e_{min}$ .

estimated leaves with  $N^l$  labeled leaves, a  $N^e \times N^l$  matrix  $D$  is computed, which records all tip-based errors of each estimated leaf tips  $\hat{t}_{1,2}$  with every labeled tips  $t_{1,2}$  normalized by the labeled leaf length:

$$e_{la}(\hat{t}_{1,2}, t_{1,2}) = \frac{\|\hat{t}_1 - t_1\|_2 + \|\hat{t}_2 - t_2\|_2}{2\|\hat{t}_1 - \hat{t}_2\|_2}. \quad (16)$$

**Algorithm 2:** Performance evaluation process.

**Input:** Tracking results  $\hat{T}$ , label results  $T$ .

**Output:**  $F$ ,  $E$ , and  $T$ .

Initialize  $f = 0$ ,  $e_1 = e_2 = []$ .

**1. for**  $i = 1, \dots, m$  **do**

$ER = \text{cell}(N^e, N^l)$ ,  $ID = \mathbf{0}_{N^e \times N^l}$ .

**for**  $j = 1, \dots, n$  **do**

$\hat{T} = \hat{T}\{i, j\}$ ;  $T = T\{i, j\}$ ;

$[f_0, ER_0, ID_0] = \text{leafMatch}(\hat{T}, T)$ ;

$f = f + f_0$ ;  $e_1 = [e_1, ER_0(ER_0 \neq 0)]$ ;

$ER = ER + ER_0$ ;  $ID = ID + ID_0$ ;

**for**  $k = 1, \dots, \min(N^e, N^l)$  **do**

$[ID_{max}, i, j] = \max(ID)$ ;

$ID(i, :) = \text{Inf}$ ;  $ID(:, j) = \text{Inf}$ ;

$e_2 = [e_2, ER\{i, j\}]$ ;

**2. for**  $\tau = 0 : 0.01 : 1$  **do**

$F(\tau) = \frac{f + \text{sum}(e_1 > \tau)}{N^b}$ ;  $E(\tau) = \text{mean}(e_1 \leq \tau)$ ;

$T(\tau) = \frac{\text{sum}(e_2 \leq \tau)}{N^b}$ .

We build the leaf correspondence by finding a number of  $\min(N^e, N^l)$  minimum errors in  $D$  that do not share columns or rows, which results in  $\min(N^e, N^l)$  leaf pairs and  $f = |N^l - N^e|$  leaves without correspondence. Finally, it outputs the number of unmatched leaf  $f$ ,  $ER$  recording tip-based errors and  $ID$  recording the leaf correspondence. This frame-to-frame correspondence is built on all 3 frames and the results are added into  $ER$  and  $ID$ . We build the video-to-video leaf correspondence using the accumulated  $ID$ .  $e_1$  and  $e_2$  are the tip-based errors of leaf pairs with frame-to-frame and video-to-video correspondence respectively. The difference of  $e_1$  and  $e_2$  is from estimated leaf 4. While it is well aligned with labeled leaf 1 in frame 3, it does not have leaf correspondence in all 3 frames together.

Finally we compute three metrics by varying a threshold  $\tau$ . **Unmatched leaf rate**  $F$  is the percentage of unmatched leaves w.r.t. the total number of labeled leaves  $N^b$ .  $F$  attributes to two sources,  $f$  leaves without correspondence and correspondent leaves with tip-based errors larger than  $\tau$ . **Landmark error**  $E$  is the average of all tip-based errors in  $e_1$  that are smaller than  $\tau$ . **Tracking consistency**  $T$  is the percentage of leaf pairs whose tip-based errors in  $e_2$  are smaller than  $\tau$  w.r.t.  $N^b$ . These three metrics jointly estimate the accuracy in leaf counting ( $F$ ), alignment ( $E$ ), and tracking ( $T$ ).

In order to quantitatively evaluate the segmentation accuracy, we annotate each image to generate a leaf segmentation mask where the pixels of the same leaf are assigned with the same number over the video. We add the metric "Symmetric Best Dice" (SBD) [5] to compute the similarity between the estimated and the ground truth segmentation masks. It is averaged across all labeled frames. These four metrics are used to evaluate the performance of our joint framework.



## 5 EXPERIMENTS AND RESULTS

### 5.1 Dataset and Templates

Our dataset includes 41 *Arabidopsis Thaliana* videos taken in a 5-day period, which is sufficient to model the plant growth [42]. Each video has 389 frames, with the image resolution ranging from  $40 \times 40$  to  $100 \times 100$ . For each video, we label the two tips of all visible leaves and segmentation masks of 5 frames, each being the middle frame of a day. In total we labeled  $N^b = 1,807$  leaves. We select 10 videos to form the training set for template generation and parameter tuning. The remaining 31 videos are used for testing. The collection of all labeled tips and segmentation masks are denoted as  $\mathbb{T}$  and  $\mathbb{B}$ .

To generate leaf templates, we select 10 leaves with representative shapes and label the two tips for each leaf, as in Figure 4. We select 12 scales for each leaf shape to guarantee the scaled templates can cover all possible leaf sizes in the dataset. For each scaled leaf template, we rotate it every  $15^\circ$  in the  $360^\circ$  space. Finally, the total number of leaf templates is  $N = SHR = 2,880$  with  $S = 10, H = 12, R = 24$ <sup>1</sup>.

### 5.2 Experimental Setup

For each testing video, we apply our approach and compare with four methods: *Baseline Chamfer Matching*, *Prior Work* [12], [16], and *Manual Results*.

**Proposed Method**  $N$  templates are applied to the edge map of the last video frame to generate the same amount of transformed templates. Leaf segmentation and alignment generate  $N^e$  leaf candidates for leaf tracking, which iteratively updates  $\mathcal{P}$  according to Equation 11 towards the first frame.

**Baseline Chamfer Matching** The basic idea of CM is to align one object in an image. To align multiple leaves in a plant image, we design the *baseline CM* to iteratively align one leaf at one time. In each iteration, we apply all  $N$  templates to the edge map of a test image to generate  $N$  transformed leaf templates, which is the same as our first step. The transformed template with the minimum CM distance is selected and denoted as a leaf candidate. We update the edge map by deleting the matched edge points of the selected leaf candidate. The iteration continues until 90% of the edge points are deleted. We apply this method to the labeled frames of each video and build the leaf correspondence based on leaf centers.

**Multi-leaf Alignment** [12] The optimization in [12] is the same as our proposed leaf alignment on the last frame. We apply [12] on all labeled frames and build the leaf correspondence based on leaf center distances.

**Multi-leaf Tracking** [16] The difference between the proposed method and [16] includes the modified  $G_3$

in Equation 11. And [16] do not have the scheme to generate a new leaf candidate during tracking.

**Manual Results** In order to find the upper bound of our proposed method, we use the ground truth labels  $\mathbb{T}$  to find the optimal set of  $\hat{\mathbb{T}}$ . For each labeled leaf, we find the leaf candidate with the smallest tip-based error  $e_{la}$  from  $N$  transformed templates.

For all methods, we record the estimated tip coordinates of all leaf candidates in the labeled frames as  $\hat{\mathbb{T}}$ . The transformed template masks are used to generate an estimated segmentation mask for each frame. We record the estimated segmentation masks of all labeled frames as  $\hat{\mathbb{B}}$ .  $\hat{\mathbb{T}}$  and  $\mathbb{T}$  are used to evaluate  $F$ ,  $E$ , and  $T$ .  $\hat{\mathbb{B}}$  and  $\mathbb{B}$  are used to evaluate  $SBD$ .

### 5.3 Experimental Results

#### 5.3.1 Performance comparison

**Qualitative Results** Figure 8 shows the results on the labeled frames of one video. Since the baseline CM only considers CM distance to segment each leaf separately, leaf candidates are likely to be aligned around the edge points, which result in large landmark errors. While [16] can keep the leaf ID consistent, it does not include the scheme to generate a new leaf candidate during tracking (e.g., leaf 8 in Figure 8). Our proposed method performs substantially better than others. It has the same segmentation as the labeled results and all leaves are well tracked. Leaf 5 is deleted when it gets too small. Due to the limitation of a *finite* amount of templates, the manual results are not perfect. However, in our tracking method, we allow template transformation under any parameters in  $\mathcal{P}$  without limiting to a finite number.

**Quantitative Results** We first evaluate the SAT accuracy w.r.t.  $F$ ,  $E$ , and  $T$ . We set the threshold  $\tau$  to vary in  $[0.05 : 0.01 : 1]$  and generate the accuracy curves for all methods, as shown in Figure 9. When  $\tau$  is small, *i.e.*, we have very strict requirements on the accuracy of tip estimation, all methods work well for easy-to-align leaves. With the increase of  $\tau$ , more and more hard-to-align leaves with relatively large tip-based errors are considered as well-aligned leaves and contribute to the landmark error  $E$  and tracking consistency  $T$ . Therefore, detecting more leaves will result in higher  $E$  and  $T$ . It is noteworthy that our method achieves *lower* landmark error and *higher* tracking consistency while segmenting *more* leaves.

The baseline CM segments less leaves with higher landmark error and lower tracking consistency. The manual results are the upper bound of our algorithm. Obviously  $F$  will be 0 and  $T$  will be 1 with the increase of  $\tau$  because we enforce the correspondence of all labeled leaves. But  $E$  will not be 0 due to the limitation of a finite template set. Overall, the proposed method performs much better than the baseline CM and our prior work. The improvement over [12] is mainly in a higher  $T$ , and it improves [16] in all

1. The dataset, labels, and templates are publicly available at: <http://cvlab.cse.msu.edu/project-plant-vision.html>.

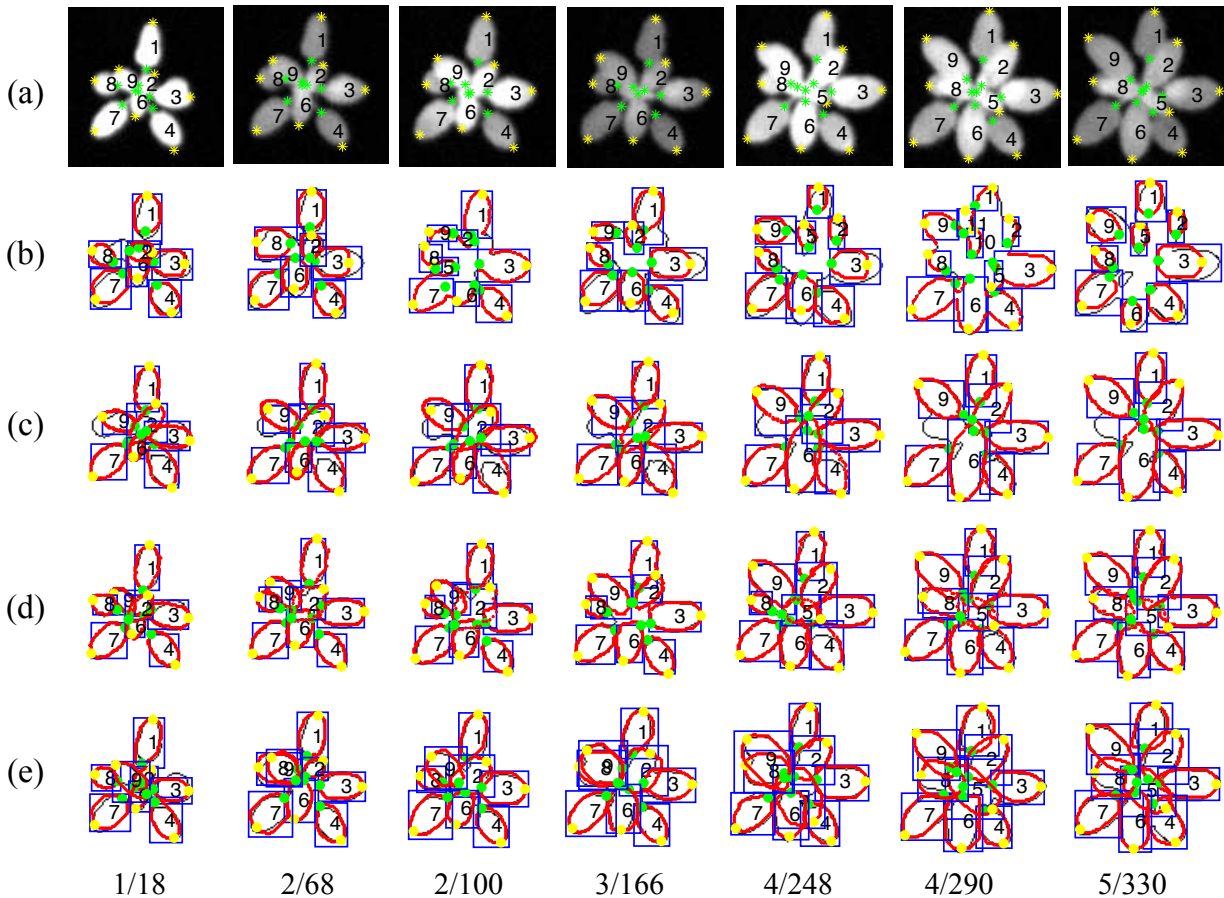


Fig. 8. Qualitative results: (a) ground truth labels; (b) baseline CM; (c) [16]; (d) proposed method; and (e) manual results. Each column is one frame in the video (day/frame). Yellow/green dots are the estimated outer/inner leaf tips. Red contour is  $W(U; p)$ . Blue box encloses the edge points matching  $W(U; p)$ . The number on a leaf is the leaf ID. Best viewed in color.

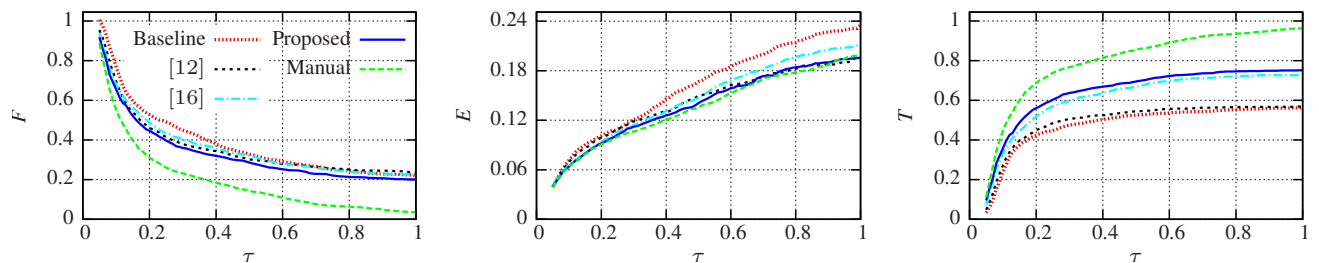


Fig. 9. Accuracy comparison of  $F$ ,  $E$ , and  $T$  vs.  $\tau$  for all different methods based on Algorithm 2.

three metrics. However there is still a gap between the proposed method and the manual results, which calls for future research.

The  $SBD$ -based segmentation accuracy is shown in Table 2. The proposed method is again superior to the baseline algorithm and the prior work.

**Efficiency Results** Table 2 shows the average execution time, which is calculated based on a Matlab implementation on a conventional computer. Our method is superior to the baseline CM and [12]. It is a little slower than [16] because of the updated  $G_3$  and the scheme to add leaf candidates during tracking.

**Segmentation Accuracy** While there is no prior work focuses on the joint multi-leaf SAT, leaf segmentation

TABLE 2  
SBD and efficiency comparison (sec./image).

	Baseline	[12]	[16]	Proposed	Manual
SBD	61.0	63.0	64.4	65.2	74.9
Time	51.28	16.42	1.98	2.15	-

TABLE 3  
Leaf segmentation SBD accuracy ( $\pm$ std) comparison.

	A1	A2	A3	all
[17]	74.2( $\pm$ 7.7)	80.6( $\pm$ 8.7)	61.8( $\pm$ 19.1)	73.5( $\pm$ 11.5)
Ours	78.5( $\pm$ 5.5)	77.4( $\pm$ 8.1)	76.1( $\pm$ 14.1)	78.0( $\pm$ 7.8)

has been studied especially on RGB imagery. For example, state-of-the-art performance [17] is reported

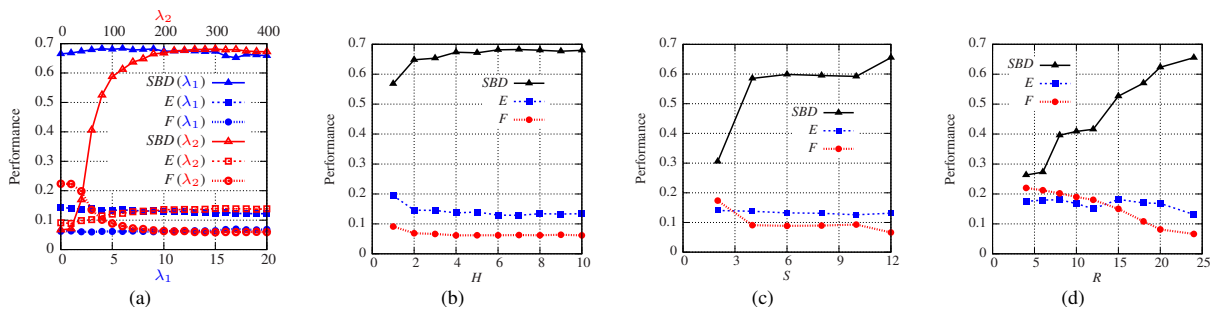


Fig. 10. Alignment parameter tuning. We show the accuracy in  $E$ ,  $F$ , and  $SBD$  when varying the coefficients of each objective (a) and the template set size (b,c,d).

in the 2014 Leaf Segmentation Challenge (LSC) [30]. We apply our segmentation and alignment algorithm to the LSC dataset [5], which consists of three sets of *Arabidopsis* (A1, A2) and *tobacco* (A3). Two examples from A2 and A3 are shown in Figure 11. Note that pre-processing of the RGB imagery is employed in order to extend our proposed method to this LSC dataset. We compare the segmentation accuracy with [17] in Table 3. Our algorithm achieves higher  $SBD$  in A1, A3, and in average. The segmentation accuracy on the LSC dataset is much higher than that of our fluorescence dataset because images in [5] are of higher resolution.

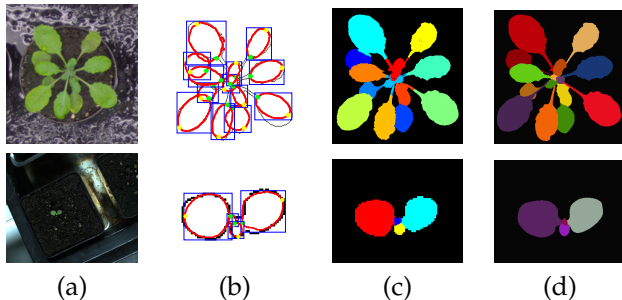


Fig. 11. Leaf segmentation results on LSC: (a) input image; (b) alignment result; (c) estimated segmentation mask; (d) ground truth segmentation mask.

### 5.3.2 Parameter tuning

We explore the sensitivity of the parameters in our method. We use the 10 training videos for parameter tuning in our framework. For alignment parameter tuning, we test on all labeled frames independently and evaluate the accuracy without using tracking consistency  $T$ . For tracking parameter tuning, we test on the labeled frames of each video and evaluate the accuracy using all four metrics.

Figure 10 (a) shows the alignment parameter tuning results of the weights for each objective term in Equation 6. We first search for the optimal setting to be:  $\lambda_1 = 4$  and  $\lambda_2 = 300$ . We then fix one parameter and change the other and evaluate the performance at  $\tau = 0.4$ . We observe that  $\lambda_1$  is relatively robust with some improvement from 0 to 4. The performance increases tremendously as  $\lambda_2$  increases, indicating that

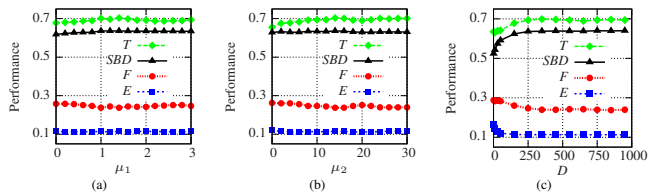


Fig. 12. Tracking parameter tuning. Accuracy w.r.t. the coefficients in the objective and the number of iterations.

$J_3$  is crucial. Without either term ( $\lambda_1 = 0$  or  $\lambda_2 = 0$ ), the performance is not optimal.

In order to analyze the impact of the number of leaf templates, we reduce the value in one of  $H$ ,  $S$ , and  $R$  at a time. As shown in Figure 10, the performance increases as the number of templates increases in all three parameters. However, orientation is the most important as leaves with different orientations are more likely to have higher CM distances than leaves with different shapes or scales.

Figure 12 (a,b) shows the parameter tuning results in leaf tracking framework. Similarly, we first find the optimal weights to be:  $\mu_1 = 1$  and  $\mu_2 = 15$ . We fix one parameter and change the other and evaluate the performance at  $\tau = 0.4$ .  $\mu_1$  and  $\mu_2$  are relatively robust to changes. However, they are still useful as without either term, the performance decreases.

To study the impact of the number of iterations between two frames, we change  $D$  and evaluate the performance. As shown in Figure 12 (c), the performance increases as  $D$  increases. However, it stabilizes when  $D$  is larger than 300 because the algorithm already converges before reaching the maximum iteration.

In summary, all parameters used in our algorithm are set as:  $\lambda_1 = 4$ ,  $\lambda_2 = 300$ ,  $C = 3$ ,  $\mu_1 = 1$ ,  $\mu_2 = 15$ ,  $D = 300$ ,  $\alpha_1 = 0.001$ ,  $\alpha_2 = 0.001$ , and  $s = 40$ .

### 5.3.3 Quality prediction

**Alignment Quality Model** Data samples for evaluating our alignment quality model are selected from  $e_1$  in Algorithm 2, which contains the tip-based errors of all leaf pairs with 90% of them are less than 0.5. We select 100 samples from  $e_1$  for each interval of tip-based error within  $[0 : 0.1 : 0.5]$ . Sample duplication is employed when the number of sample in a particular interval is less than 100. All samples with tip-based error larger than 0.5 will also be selected but without

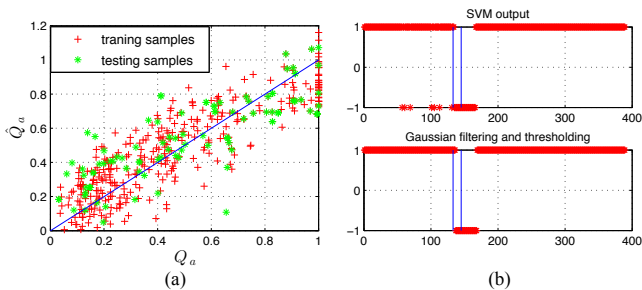


Fig. 13. (a) Alignment quality model applied to training and testing samples; (b) Tracking quality model applied to one video: the SVM classifier output (top), the result after Gaussian filtering and thresholding (bottom), and the two blue lines are the labeled period of a tracking failure.

duplication. Finally we select 625 samples and extract features  $x_a$  for each sample. We assign  $Q_a = 2e_{la}$  to make the output in the range of  $[0, 1]$ . And  $Q_a = 1$  for all samples with  $e_{la} > 0.5$ . We randomly select 100 samples as the test set and the remaining samples are used to train the model. Figure 13 (a) shows the results of the model on both training and testing samples.

We use  $R^2$  to measure how well the model fits our data. It is defined as:

$$R^2 = 1 - \frac{\sum (Q_a - \hat{Q}_a)^2}{\sum (Q_a - \bar{Q}_a)^2}, \quad (17)$$

where  $\hat{Q}_a$  is the predicted quality value and  $\bar{Q}_a$  is the mean of  $Q_a$ . In our model,  $R^2 = 0.769$  and the correlation coefficients for all testing samples is 0.813. Both values indicate a high correlation of  $Q_a$  and  $\hat{Q}_a$ . This quality model is used to predict the alignment accuracy and generate one predicted curve for each leaf, as shown in Figure 2.

**Tracking Quality Model** We visualize the results of our method and find 15 videos that have a tracking failure of one leaf. As the goal for tracking quality model is to detect when the tracking failure starts, we label two frames when the failure starts and ends in each video. The starting frame is when a leaf candidate starts to change its location toward its neighboring leaves. The ending frame is when a leaf candidate totally overlaps another leaf. Among all failure samples, the shortest tracking failure length is 5 frames and the average length is 12 frames.

We select 2-3 frames near the ending frame as the negative training samples with  $Q_t = -1$  and 5 frames evenly distributed before the failure starts as the positive training samples with  $Q_t = 1$ . The features  $x_t$  are extracted as discussed in Section 3.3.2 and used to train a SVM classifier. The learned model is applied to all frames to predict the tracking quality. Figure 13 (b) shows an example of the output. We apply a Gaussian filter to remove outliers and delete the failure whose length is less than 5 frames (the shortest length of failure samples).

We compare the first frame of a predicted failure

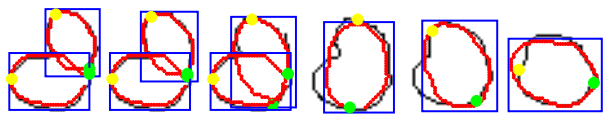


Fig. 14. Leaf alignment results on synthetic leaves with various amount of overlap. From left to right, the overlap ratio w.r.t. the smaller leaf is 10%, 15%, 22%, 23%, 36%, and 59%.

with that of a labeled failure. When their distance is less than 12 frames (the average length of the failure samples), it is considered as a true detection. Otherwise it is a false detection. Using the leave-one-video-out testing scheme, the quality model generates 11 true detections and 16 false detections over 15 labeled failures. Similarly, this quality model is applied during tracking and outputs a prediction curve for each leaf (shown in Figure 2).

### 5.3.4 Limitation analysis

Any algorithm has its limitation. Hence, it is important to explore the limitation of the proposed method. First, one interesting question is to what extend our segmentation and alignment method can correctly segment leaves in the overlapping region. We answer this question using a simple synthetic example. As shown in Figure 14, our method performs well when the overlap ratio is less than 23%. Otherwise it identifies two leaves as one leaf, which appears to be reasonable when the overlap ratio is high (e.g., 59%).

Second, our leaf tracking starts from a good initialization of the leaf candidates from the previous frame. Another interesting question is to what extend our tracking method can succeed with bad initializations. To study this, 6 frames with good tracking results are selected from 6 videos (one for each). We change the transformation parameters  $P$  to synthesize different amount of distortions and apply our tracking algorithm on these 6 frames. The leaf candidate is deleted only if it becomes one point and the tip-based error is set to be 1. We compute the average tip-based error of all leaf candidates.

We vary the rotation angle  $\theta$ , the scaling factor  $r$ , and the translation ratio  $t_{xy}$ , which is defined as  $t_{xy} = \frac{\sqrt{t_x^2 + t_y^2}}{\sqrt{(t_1^x - t_2^x)^2 + (t_1^y - t_2^y)^2}}$  and the direction is randomly selected. The average and range of the tip-based errors for all 6 frames are shown in Figure 15. Our tracking method reduces the initial tip-based error to a small value. It is most robust to  $r$  and most sensitive to  $t_{xy}$ .

Figure 16 shows some examples. For rotation angle less than  $45^\circ$ , our method works well for different amounts of leaf rotations. For the scaling factor, as long as the leaf candidate is not too small, our method is very robust even if we enlarge the original leaf candidates to be 2.5 times larger. For the translation ratio, it is sensitive because the direction is randomly selected and leaf candidates are very likely to shift to

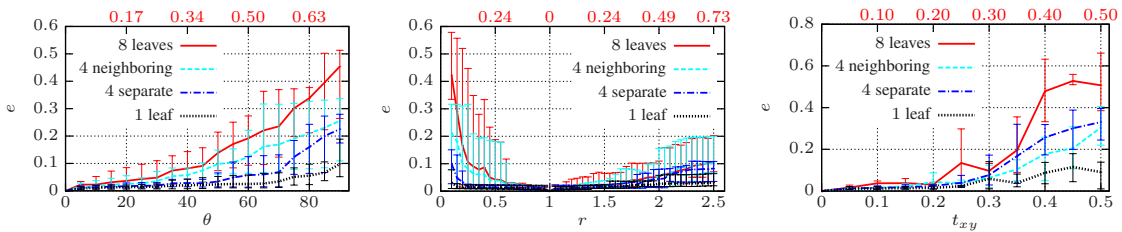


Fig. 15. Mean tip-based error with different initializations. The axes on top of the figures show the initial tip-based errors.

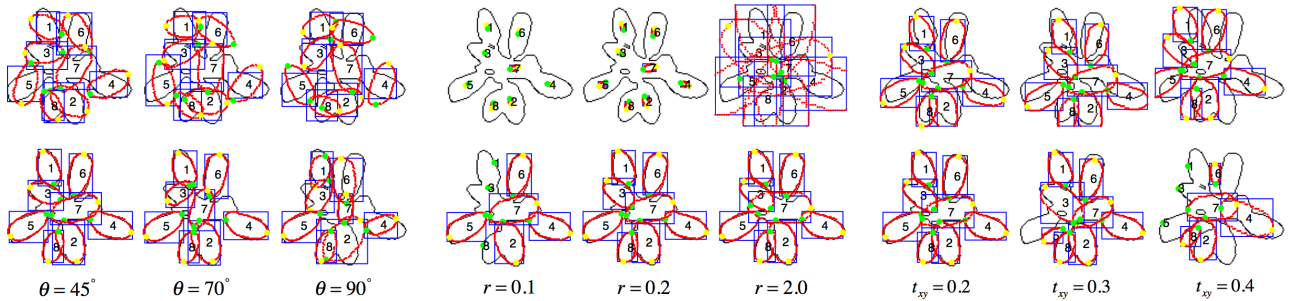


Fig. 16. Example results: the first row shows the initialization, and the second row shows the tracking results.

the locations of the neighboring leaves. Furthermore, changing the initialization of  $\theta$  and  $r$  for 4 separate leaves (leaf 1, 4, 5, 8 in Figure 16) leads to better performance than that of 4 neighboring leaves (leaf 1, 3, 6, 7 in Figure 16) because neighboring leaves will have overlap with each other and therefore influence the tracking results. Overall, as the distortion increases, the average tip-based error increases while some of the leaf candidates can still be well aligned.

## 6 CONCLUSIONS

In this paper, we identify a new computer vision problem of leaf segmentation, alignment, and tracking from fluorescence plant videos. Leaf alignment and tracking are formulated as two optimization problems based on Chamfer matching and leaf template transformation. Two models are learned to predict the quality of leaf alignment and tracking. A quantitative evaluation scheme is designed to evaluate the performance. The limitations of our algorithm are studied and experimental results show the effectiveness, efficiency, and robustness of the proposed method.

With the leaf boundary and structure information over time, the *photosynthetic efficiency* can be computed for each leaf, which paves the way for leaf-level photosynthetic analysis and high-throughput plant phenotyping. The proposed method and the evaluation scheme are potentially applicable to other plant videos, as shown in the results on the LSC dataset.

## REFERENCES

- [1] Fabio Fiorani and Ulrich Schurr, "Future scenarios for plant phenotyping," *Annual Review of Plant Biology*.
- [2] Marcus Jansen et al., "Simultaneous phenotyping of leaf growth and chlorophyll fluorescence via growSCREEN fluoro allows detection of stress tolerance in *Arabidopsis thaliana* and other rosette plants," *Functional Plant Biology*, 2009.
- [3] Samuel Trachsel, Shawn M Kaeppeler, Kathleen M Brown, and Jonathan P Lynch, "Shovelomics: high throughput phenotyping of maize (*Zea mays* L.) root architecture in the field," *Plant and Soil*, 2011.
- [4] Larissa M Wilson, Sherry R Whitt, Ana M Ibáñez, Torbert R Rocheford, Major M Goodman, and Edward S Buckler, "Dissection of maize kernel composition and starch production by candidate gene association," *The Plant Cell*, 2004.
- [5] Hanno Scharf, Massimo Minervini, Andreas Fischbach, and Sotirios A Tsafaris, "Annotated image datasets of rosette plants," Tech. Rep. FZJ-2014-03837, 2014.
- [6] Anja Hartmann, Tobias Czauderna, Roberto Hoffmann, Nils Stein, and Falk Schreiber, "Htpheno: an image analysis pipeline for high-throughput plant phenotyping," *BMC Bioinformatics*, 2011.
- [7] Ladislav Nedbal and John Whitmarsh, "Chlorophyll fluorescence imaging of leaves and fruits," in *Chlorophyll a Fluorescence*. 2004.
- [8] Xu Zhang, Ronald J. Hause, and Justin O. Borevitz, "Natural genetic variation for growth and development revealed by high-throughput phenotyping in *Arabidopsis thaliana*," *G3: Genes, Genomes, Genetics*, 2012.
- [9] Arabidopsis Genome Initiative et al., "Analysis of the genome sequence of the flowering plant *Arabidopsis thaliana*," *Nature*.
- [10] "http://www.arabidopsis.org/portals/education/aboutarabidopsis.jsp#hist."
- [11] Chin-Hung Teng, Yi-Ting Kuo, and Yung-Sheng Chen, "Leaf segmentation, its 3D position estimation and leaf classification from a few images with very close viewpoints," in *Image Analysis and Recognition*. 2009.
- [12] Xi Yin, Xiaoming Liu, Jin Chen, and David M Kramer, "Multi-leaf alignment from fluorescence plant images," in *WACV*, 2014.
- [13] Jonas Vyllder, Daniel Ochoa, Wilfried Philips, Laury Chaerle, and Dominique Straeten, "Leaf segmentation and tracking using probabilistic parametric active contours," in *Computer Vision/Computer Graphics Collaboration Techniques*. 2011.
- [14] Harry G. Barrow, Jay M. Tenenbaum, Robert C. Bolles, and Helen C. Wolf, "Parametric correspondence and Chamfer matching: Two new techniques for image matching," Tech. Rep., DTIC Document, 1977.
- [15] Bastian Leibe, Edgar Seemann, and Bernt Schiele, "Pedestrian detection in crowded scenes," in *CVPR*, 2005.
- [16] Xi Yin, Xiaoming Liu, Jin Chen, and David M Kramer, "Multi-leaf tracking from fluorescence plant videos," in *ICIP*, 2014.
- [17] Jean-Michel Pape and Christian Klukas, "3-d histogram-based segmentation and leaf detection for rosette plants," in *ECCV Workshops*, 2014.

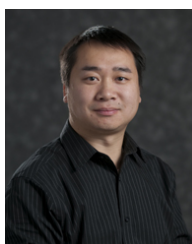
- [18] Long Quan, Ping Tan, Gang Zeng, Lu Yuan, Jingdong Wang, and Sing Bing Kang, "Image-based plant modeling," in *ACM Transactions on Graphics*, 2006.
- [19] Derek Bradley, Derek Nowrouzezahrai, and Paul Beardsley, "Image-based reconstruction and synthesis of dense foliage," *ACM Transactions on Graphics*.
- [20] Yangyan Li, Xiaochen Fan, Niloy J Mitra, Daniel Chamovitz, Daniel Cohen-Or, and Baoquan Chen, "Analyzing growing plants from 4d point cloud data," *ACM Transactions on Graphics*, 2013.
- [21] Yann Chéné, David Rousseau, Philippe Lucidarme, Jessica Bertheloot, Valérie Caffier, Philippe Morel, Étienne Belin, and François Chapeau-Blondeau, "On the use of depth camera for 3D phenotyping of entire plants," *Computers and Electronics in Agriculture*, 2012.
- [22] Guillaume Cerutti, Laure Tougne, Julien Mille, Antoine Vacavant, and Didier Coquin, "Understanding leaves in natural images—a model-based approach for tree species identification," *CVIU*, 2013.
- [23] Sofiene Mouine, Itheri Yahiaoui, and Anne Verroust-Blondet, "Advanced shape context for plant species identification using leaf image retrieval," in *Proc. ACM Int. Conf. Multimedia Retrieval (ICMR)*, 2012.
- [24] Neeraj Kumar, Peter N. Belhumeur, Arijit Biswas, David W. Jacobs, W. John Kress, Ida C. Lopez, and João VB. Soares, "Leafsnap: A computer vision system for automatic plant species identification," in *ECCV*, 2012.
- [25] Guillaume Cerutti, Laure Tougne, Julien Mille, Antoine Vacavant, Didier Coquin, et al., "A model-based approach for compound leaves understanding and identification," in *ICIP*, 2013.
- [26] Guillaume Cerutti, Laure Tougne, Antoine Vacavant, and Didier Coquin, "A parametric active polygon for leaf segmentation and shape estimation," in *Advances in Visual Computing*, 2011.
- [27] Siqi Chen, Daniel Cremers, and Richard J Radke, "Image segmentation with one shape prior template-based formulation," *Image and Vision Computing*, 2012.
- [28] Xiao-Feng Wang, De-Shuang Huang, Ji-Xiang Du, Huan Xu, and Laurent Heutte, "Classification of plant leaf images with complicated background," *Applied Mathematics and Computation*, 2008.
- [29] Xianghua Li, Hyo-Haeng Lee, and Kwang-Seok Hong, "Leaf contour extraction based on an intelligent scissor algorithm with complex background," in *International Conference on Future Computers in Education*, 2012.
- [30] "http://plant-phenotyping.org/CVPPP2014-challenge."
- [31] Hanno Schar, Massimo Minervini, Andrew P French, Christian Klukas, David M Kramer, Xiaoming Liu, Imanol Luengo, Jean-Michel Pape, Gerrit Polder, Danijela Vukadinovic, et al., "Leaf segmentation in plant phenotyping: a collation study," *Machine Vision and Applications*, 2015.
- [32] Eren Erdal Aksoy, Alexey Abramov, Florentin Wörgötter, Hanno Schar, Andreas Fischbach, and Babette Dellen, "Modeling leaf growth of rosette plants using infrared stereo image sequences," *Computers and Electronics in Agriculture*, 2015.
- [33] Babette Dellen, Hanno Schar, and Carme Torras, "Growth signatures of rosette plants from time-lapse video," 2015.
- [34] Martin A Fischler and Robert C Bolles, "Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography," *Communications of the ACM*, 1981.
- [35] Xiaoming Liu, Ting Yu, Thomas Sebastian, and Peter Tu, "Boosted deformable model for human body alignment," in *CVPR*, 2008.
- [36] Xiaoming Liu, "Discriminative face alignment," *PAMI*, 2009.
- [37] Yves Crama, Pierre Hansen, and Brigitte Jaumard, "The basic algorithm for pseudo-boolean programming revisited," *Discrete Applied Mathematics*, 1990.
- [38] Endre Boros and Peter L Hammer, "Pseudo-boolean optimization," *Discrete Applied Mathematics*, 2002.
- [39] Eyung Lim, Xudong Jiang, and Weiyun Yau, "Fingerprint quality and validity analysis," in *ICIP*, 2002.
- [40] Kamal Nasrollahi and Thomas B. Moeslund, "Face quality assessment system in video sequences," in *Biometrics and Identity Management*. 2008.
- [41] Kathleen Greenham, Ping Lou, Sara E Remsen, Hany Farid, and C Robertson McClung, "Trip: Tracking rhythms in plants, an automated leaf movement analysis program for circadian period estimation," *Plant Methods*, 2015.
- [42] Oliver L Tessmer, Yuhua Jiao, Jeffrey A Cruz, David M Kramer, and Jin Chen, "Functional approach to high-throughput plant growth analysis," *BMC systems biology*, 2013.



**Xi Yin** received the B.S. degree in Electronic and Information Science from Wuhan University, China, in 2013. Since August 2013, she has been working toward her Ph.D. degree in the Department of Computer Science and Engineering, Michigan State University, USA. Her paper on plant segmentation won the Best Student Paper Award at Winter Conference on Application of Computer Vision (WACV) 2014. Her research interests include computer vision and deep learning.



**Xiaoming Liu** is an Assistant Professor in the Department of Computer Science and Engineering at Michigan State University (MSU). He received the B.E. degree from Beijing Information Technology Institute, China and the M.E. degree from Zhejiang University, China, in 1997 and 2000 respectively, both in Computer Science, and the Ph.D. degree in Electrical and Computer Engineering from Carnegie Mellon University in 2004. Before joining MSU in Fall 2012, he was a research scientist at General Electric Global Research Center. His research areas are face recognition, biometrics, image alignment, video surveillance, computer vision and pattern recognition. He has authored more than 70 scientific publications, and has filed 22 U.S. patents. He is a member of the IEEE.



**Jin Chen** received the BS degree in computer science from Southeast University, China, in 1997, and the Ph.D. degree in computer science from the National University of Singapore, Singapore, in 2007. He is an Assistant Professor in the Department of Energy Plant Research Laboratory and the Department of Computer Science and Engineering at Michigan State University. His general research interests are in computational biology, as well as its interface with data mining and computer vision.



**David M. Kramer** is the Hannah Distinguished Professor of Bioenergetics and Photosynthesis in the Biochemistry and Molecular Biology Department and the MSU-Department of Energy Plant Research Lab at Michigan State University. In 1990, he received his Ph.D. in Biophysics at University of Illinois, Urbana-Champaign, followed by Post-doc at the Institute de Biologie Physico-Chimique in Paris and a 15-year tenure as a faculty member at Washington State University. His research seeks to understand how plants convert light energy into forms usable for life, how these processes function at both molecular and physiological levels, how they are regulated and controlled, how they define the energy budget of plants and the ecosystem and how they have adapted through evolution to support life in extreme environments. This work has led his research team to develop a series of novel spectroscopic tools for probing photosynthetic reactions in vivo.