
1 Non-destructive Automatic Leaf Area Measurements by Combining Stereo and 2 Time-of-Flight Images

3 Yu Song, Chris A. Glasbey, Gerrit Polder, Gerie W.A.M. van der Heijden

4

5 **Abstract** Leaf area measurements are commonly obtained by destructive and laborious practice. This paper
6 shows how stereo and Time-of-Flight (ToF) images can be combined for non-destructive automatic leaf area
7 measurements. We focus on some challenging plant images captured in a greenhouse environment, and show that
8 even the state-of-the-art stereo methods produce unsatisfactory results. By transforming depth information in a
9 ToF image to a localised search range for dense stereo, a global optimisation strategy is adopted for producing
10 smooth results that preserve discontinuity. We also use edges of colour and disparity images for automatic leaf
11 detection and develop a smoothing method necessary for accurately estimating surface area. In addition to
12 show that combining stereo and ToF images gives superior qualitative and quantitative results, 149 automatic
13 measurements on leaf area using our system in a validation trial have a correlation of 0.97 with true values
14 and the root-mean-square error is 10.97 cm^2 , which is 9.3% of the average leaf area.. Our approach could
15 potentially be applied for combining other modalities of images with large difference in image resolutions and
16 camera positions.

17 **Keywords** Dense Stereo · Time-of-Flight · Leaf Detection · Surface Reconstruction · 3D Measurements

18 1 Introduction

19 In our post-genomic world, where we are deluged with genetic information, the bottleneck to scientific progress
20 is often phenotyping, i.e. measuring the observable characteristics of plants and animals. For example, surface

Y. Song · C. A. Glasbey
Biostatistics and Statistics Scotland, The King's Buildings, Edinburgh, EH9 3JZ, UK
E-mail: {y.song,chris.glasbey}@bioss.ac.uk

G. Polder · G.W.A.M. van der Heijden
Biometris, Wageningen UR, PO Box 100, 6700 AC Wageningen, Netherlands
E-mail: {gerrit.polder,gerie.vanderheijden}@wur.nl

21 area of a leaf is a powerful diagnostic of plant productivity. Current common practice, also known as direct
22 measurement, requires each plant leaf to be manually stripped and fed through a dedicated measuring machine,
23 which is destructive, laborious and time-consuming [23]. Owing to these difficulties, leaf area measurements
24 during the plant growth period as in [46] are often not possible and they have been empirically simulated in
25 plant growth analysis [21,32].

26 Image analysis has the potential to overcome these problems. The aim is to develop an automated procedure
27 for extracting each leaf individually from an image, followed by reconstruction and measurement of extracted
28 leaves in 3D. Automatic interpretation of plant images remains very difficult. Fig. 1 shows a stereo pair of images
29 of pepper plants captured using the setup shown in Fig. 2, using flash lighting both to cancel the variable effects
30 of natural illumination and permit a fast shutter speed which minimises blur while the camera rig is in motion.
31 Our aim is to recover dense depth information and estimate leaf area in 3D along with other plant characteristics
32 such as stem length or fruit size. This is a challenging task, as surfaces are of complex shape, and there are
33 multiple depths, linear features, occlusions and inconsistent shadows between images. Stereo vision is usually
34 seemingly-effortless for the human eye and brain, but unfortunately still not so for computers.

35 Recently the use of low-resolution range cameras based on the Time-of-Flight (ToF) principle has received
36 increasing attention, with Kolb *et. al.* [25] providing an overview on techniques and potential applications.
37 The state-of-the-art stereo vision algorithms are known for underperforming in low-textured areas (usually at
38 object centres) and preserving edge discontinuity, which can be complemented by the use of a ToF camera.
39 Augmenting stereo pairs with ToF images can therefore be used for applications aimed of improving recovery of
40 depth information. However, as in Fig. 1, the ToF camera has different field-of-view and position to the colour
41 camera, and the difference in resolution between the colour image and the ToF image is enormous.

42 These images were collected as part of an EU-funded FP7 project, SPICY (Smart tools for Prediction and
43 Improvement of Crop Yield). The plant breeding industry has contributed greatly to the increased quality and
44 yield of plant products over recent decades. However, to sustain and accelerate this progress, the relationship
45 between genotype and phenotype needs to be better understood. For example, yield is a result of the interaction
46 of many genetic factors, and is also subject to large, extraneous variation. The approach taken in SPICY is to
47 use crop growth models to predict the phenotypic response, with genotype encapsulated in model parameters

48 [2,20]. Our task in the project is the development of image analysis tools to replace hand measurements for
49 phenotyping over a large range of genotypes in a practical environment.

50 Fig. 3 shows the main processes of our approach for leaf area estimation. The first step is to recover dense
51 depth information from image pairs, since leaves appear in different sizes in an image at different depths. Without
52 depth estimation, we can only calculate projected area of leaves as in [38]. Leaf detection enables extracting
53 each individual leaf from images, and together with depth estimates, three-dimensional surface of each leaf can
54 then be reconstructed and measured. In effect, we aim to identify individual leaves and then measure them,
55 which is analogous to the manual measurement process.

56 In Fig. 1, the ratio of pixels between colour and ToF images is 200 : 1. With the development of ToF camera
57 such as a more recent camera used in [1]¹, image registration[7] in principle could become a possible solution
58 but cluttered scenes as in Fig. 1 are difficult. In this paper, we used a rigid setup allowing projecting a partial
59 and coarse-resolution ToF image into a corresponding colour image as an aid to stereo vision.

60 Part of our work has been described in [37,42,20]. [37] first described the imaging setup we built, and [42]
61 described the general idea of our work. [20] discussed the practical imaging applications using our methods and
62 how they can be used for quantifying plant characteristics. This paper presents details on camera calibration
63 for combining stereo and ToF images, and a smoothing technique for surface reconstruction. We also describe
64 a edge-based segmentation method using the 3D geometry to extract individual leaves from images.

65 Our imaging setup and our approach in principle can be applied for combining many modalities of images
66 with large difference in image resolutions and camera positions. For example, our imaging setup could have a
67 thermal camera providing temperature information, and high dimensional data from multiple imaging modalities
68 would lead to many practical applications.

69 **2 Relevant work**

70 Dense stereo produces a depth estimate for every pixel in an image, which is required in the first step of our
71 system. One approach to dense stereo is via local descriptors such as SIFT [31], followed by methods such as
72 DAISY [44] and SIFTflow [29], both of which deal with challenging stereo images. Discontinuity preserving
73 results are highly desirable for our application, but no quantitative result addressing this issue was produced in

¹ In [1], the resolutions were 640×480 for colour images and 204×204 for ToF images, and their ratio of pixels is 7.38 : 1.

74 [29,44]. Though global optimisation methods such as graph cuts [5] can produce edge-preserving results on the
75 Middlebury stereo dataset[39], challenges in the Middlebury stereo dataset are different to these images in our
76 work. Ogale and Aloimonos [34] proposed to use shape in establishing edge-preserving dense correspondence,
77 but surfaces in our images are more complicated than theirs.

78 Gudmundsson *et. al.* [17] transformed ToF points into colour images by rectification homographies, and then
79 fed them into a hierarchical stereo matching algorithm. Hahne and Alexa [18] demonstrated the combined ToF
80 and stereo method can enhance the depth estimation even without accurate extrinsic calibration. Zhu *et. al.* [47]
81 developed a weighting method combining stereo and ToF data by fixed values, and then used belief propagation
82 to optimise the data. Motivated by this research, we first present a geometric approach to transform points
83 from ToF image coordinates to colour image coordinates, and then derive a localised search range for stereo
84 matching. Despite the simplicity of the ToF transformation, we demonstrate that a global stereo strategy can
85 then be applied and does improve results and preserve discontinuity. Compared with above works [17,18,47],
86 difficult low-resolution ToF images 48×64 were used in this work. Beder *et. al.* [4] also developed a fusion
87 scheme using ToF images in the same resolution as ours except that their images were planar surfaces.

88 Current ToF and stereo fusion work (e.g. [4,17,18,27]) lack quantitative results on preserving depth discon-
89 tinuity, with the exception of [47], which used another 3D scanner to produce pixel-by-pixel depth data in an
90 indoor lab environment. In our greenhouse setting it is problematic to collect accurate pixel-by-pixel depth
91 data, so we use an indirect method to evaluate the quality of depth estimation for competing methods, by
92 quantifying how well depth-discontinuity is preserved.

93 Foreground extraction of live video using ToF and colour cameras is proposed by Wang *et. al.* [45] in order to
94 segment a person in foreground from the background. Their challenges are to track and segment the foreground
95 person from a continuous video sequence, while our images have very limited views of foreground objects. Similar
96 applications in this area have also been investigated in [8,16,40].

97 Regarding leaf measurement, existing work [35,38,43] focuses on collecting images of individual plants that
98 are separately transported from the greenhouse to a controlled imaging environment, or on imaging single leaves
99 against a plain background [26]. Instead of moving individual plants around the greenhouse, our methodology
100 brings the imaging equipment to the plants. Transporting growing plants is undesirable, because of potential
101 damage to plants that can highly disturb their growth. But more importantly, for many greenhouse crops like

102 pepper and tomato, the plants are simply too large to be transported. Our system measures plants in their own
 103 growing environment, and does not require transporting plants. However, by using our system, challenges arise
 104 from less controllable lighting conditions and a cluttered scene with large occlusions as shown in Fig. 1 and 2.

105 3 Setup and Calibration

106 Every pepper plant has a QR barcode for relating the manual measurements to the automatic measurements,
 107 and the plants grow in rows with heating pipes in-between. The maximum height of the plants is about three
 108 metres, while the space between rows is only one metre. Four camera rigs were therefore vertically stacked in a
 109 trolley known as Spy-See to cover the complete field of view, and [37, 20] provided hardware details.

110 Our camera rig consists of a colour camera and a ToF camera. The ToF camera is a radio-frequency modu-
 111 lated camera with phase shift detectors (IFM O3D201 PMD camera), with a resolution of 64×48 pixels, while
 112 the colour camera has a resolution of 480×1280 . The Spy-See setup moves in a straight line on top of rigid
 113 heating pipes in the greenhouse and captures overlapping images at a fixed interval (see Fig. 2). The baseline
 114 between images is 5 cm, and objects of interest (e.g. leaves) are located between 55 cm and 120 cm away from
 115 the camera.

116 Once assembled, our imaging setup was rigid and fixed. The positions of the colour and ToF cameras were
 117 unchanged, so was the capture interval. We therefore performed depth calibration to find both cameras in 3D
 118 space relative to each other.

119 A two-layer board shown in Fig. 4(a) was used for calibration of the colour camera at different distances
 120 from the camera. The front layer moved from 40 cm to 120 cm away from the camera in 5 cm steps. We used
 121 a simple pinhole camera model for the colour camera [19] as shown in Fig. 4(b). Let (x, y) be the position of a
 122 point in colour camera coordinate, and $(x_i, y_i), i = 0, 1, 2 \dots$ represents the position in image i . Given a point
 123 in the world coordinate (X, Y, Z) , (x_i, y_i) can be obtained as,

$$x_i = (X - s_i) f / Z + x_m \tag{1}$$

$$y_i = Y f / Z + y_m \tag{2}$$

124 where s is the baseline between the two images, which was 5 cm for our setup, and f is the focal point of
 125 the colour camera. We set the principal point (x_m, y_m) as the image centre for simplicity, which only produced
 126 negligible errors. We only consider horizontal disparity since our imaging setup moved in the horizontal direction
 127 only. Given a point identified in two colour images x_0 and x_1 , from (1),

$$x_0 = X f / Z + x_m \quad (3)$$

$$x_1 = (X - s) f / Z + x_m \quad (4)$$

128 Let $d = x_0 - x_1$ be the disparity,

$$d = s f / Z \quad (5)$$

129 During calibration, multiple depth measurements \mathbf{Z} (e.g. 40 cm to 120 cm in this work) and correspondences
 130 in each view \mathbf{d} are used to compute \hat{f} by applying the least squares fitting technique.

$$\hat{f} = \arg \min_f \| \mathbf{d} - s f / \mathbf{Z} \|^2 \quad (6)$$

131 The centre of the square seen in Fig. 4(a) is used to compute \mathbf{d} , and \mathbf{Z} is known for each image. The squares
 132 can be identified by linear discriminant analysis and connected-component labelling [13]. We manually labelled
 133 squares in two board images as training data, and linear discriminant analysis was subsequently used to segment
 134 squares for all the other board images. Since the centre of the square was required, morphological dilation was
 135 applied for refining the square shape before connected-component labelling. Fig. 4(c) presents the relationship
 136 in (5) and plots \mathbf{d} against \mathbf{Z} . Compared to a flat checkerboard used in [28], the two-layer board is also used for
 137 the ToF camera, to convert ToF measurements (x', y', z') to world coordinate (X, Y, Z) .

138 Since the relative positions of colour and ToF cameras did not change, a transformation can be established
 139 for points in colour image and ToF image. A ToF image maps Z on z' and a near-linear relationship between
 140 ToF depth measurements z' and Z was observed as in [47]. As in (1) (2), the focal point f' is still needed for
 141 the transformation of $X \rightarrow x'$ and $Y \rightarrow y'$.

$$x'_i = (X - X_0 - s i) f' / Z + x'_m \quad (7)$$

$$y'_i = (Y - Y_0) f' / Z + y'_m \quad (8)$$

142 X_0 and Y_0 are the physical distance in cm between the colour camera and the ToF camera. Although the rela-
 143 tionship between the colour camera and the ToF camera is translational in this work, [19] provided information
 144 on the homogeneous affine transformation between two cameras.

145 During ToF camera calibration, we used known Z as a cue to perform thresholding for identifying the centre
 146 of each square, and then compute d' using two adjacent ToF images. The same procedure for \hat{f} as in (6) was used
 147 to obtain \hat{f}' for the ToF camera. Challenges and error sources in ToF camera calibration have been discussed
 148 in [25, 28] ([28] also provided calibration software), and readers can review these works for further information.

149 4 Combining stereo and Time-of-Flight images

150 4.1 Dense stereo methods

151 Dense stereo methods can estimate disparity d for every pixel given a pair of stereo images. However, the pixel
 152 consistency assumption is often made for building the correspondence between two images. In our application,
 153 we have found that pixel values were not reliable for matching due to changes of perspective, lighting, and noise.
 154 To address this issue, the SIFTflow [29] method was chosen, which uses pixel-wise SIFT features between two
 155 images instead of pixel values for matching. Complex image pairs across different scenes and object appearances
 156 have been shown robustly matched in [29].

157 The pepper plant images shown in Fig. 1 have very sharp depth edges, and we have observed step changes
 158 over 50 pixels between neighbourhood pixels. Although Liu *et. al.* used a simple synthetic image in [29] to
 159 demonstrate that the dense SIFT features contain sharp edges with respect to the sharp edges in the original
 160 image, there is no close-up on complex scenes to prove that the SIFTflow method can preserve discontinuity.
 161 Ogale and Aloimonos [34] examined the implications of shape on the process of finding dense correspondence,
 162 and attempted to produce disparities in the form of a piecewise continuous function consistent with the stereo
 163 images. Using piecewise constant and piecewise linear shape models, they have presented results on images with

164 slanted planar surfaces as well as a pair of stereo images on some branches of a tree, but no results on curved
 165 or nonrigid surfaces common in the pepper plant images have been shown.

166 Global optimisation methods such as graph cuts and belief propagation have been shown producing satisfac-
 167 tory discontinuity-preserving results on the Middlebury stereo dataset [39]. Since global stereo methods produce
 168 better results compared with local stereo methods for combining with ToF information [47], we chose the alpha
 169 expansion technique applied in a graph-based energy minimisation framework [5]. The energy cost E given a
 170 pixel disparity d is defined as:

$$E(d) = \sum D(d_{(x,y)}) + \sum_{q \in N} V(d_{(x,y)}, d_{(x_q, y_q)}) \quad (9)$$

171 where N denotes the first-order neighbourhood pixels. For the data term cost D ,

$$D(d_{(x,y)}) = \min \left\{ \frac{1}{3} \sum_{c=\{R,G,B\}} \left| I_{(x,y)}^{(c)} - I'_{(x+d_{(x,y)}, y)}^{(c)} \right|, T_d \right\} \quad (10)$$

172 where I and I' represent the intensity value in the pair of colour images. T_d is a truncation constant, and
 173 $D(d_{(x,y)})$ is computed for all the possible disparities. For the smoothness term cost V ,

$$V(d_{(x,y)}, d_{(x_q, y_q)}) = u_{(x,y,x_q, y_q)} \min \left\{ |d_{(x,y)} - d_{(x_q, y_q)}|, T_k \right\} \quad (11)$$

174 where parameter T_k is used to truncate the linear energy. (x_q, y_q) is one of the first-order 4-neighbourhood pixels
 175 around (x, y) . $u_{(x,y,x_q, y_q)}$ represents static cues in Boykov *et. al.* [5], which was used as an indicator function in
 176 this work as:

$$u_{(x,y,x_q, y_q)} = \begin{cases} \alpha_v & \text{if } \sum_{c=\{R,G,B\}} \left| I_{(x,y)}^{(c)} - I_{(x_q, y_q)}^{(c)} \right| > T_e \\ n_v \alpha_v & \text{otherwise} \end{cases} \quad (12)$$

177 α_v is the smoothness cost for intensity edges produced by the thresholding value T_e , and $n_v \alpha_v$ is the smoothness
 178 cost for surfaces. Both α_v and $n_v \alpha_v$ should be set according to the data cost values in (10). (12) gives more
 179 smoothness if there is no intensity edge, and therefore achieves edge-preservation by encouraging changes at
 180 edges at a cost of α_v and limiting changes on the surface by $n_v \alpha_v$.

181 4.2 Localised search range from ToF image

182 Given the complexity associated with the pepper plant images for dense stereo methods, a localised search range
 183 $[d_{min}, d_{max}]$ derived from the corresponding ToF depth image should improve the estimation accuracy.

184 Since the ToF image is much coarser in resolution compared to the colour image, the transformation from
 185 ToF image coordinates to colour image coordinates alone would only give isolated point depth measurements in
 186 the colour image. We therefore treat each ToF pixel as a patch centring around the pixel, and then transform
 187 all points in the patch to the colour image (see ToF in Fig. 9 for an example). In effect, this transformation is
 188 one of the up-scaling techniques as discussed by Lindner *et. al.* [27] and they provided a biquadratic scheme for
 189 this purpose.

190 Due to different viewing positions of ToF and RGB cameras, there are n ToF measurements for Z ($n \geq 0$)
 191 at location (x, y) . If multiple depths were found at (x, y) , the minimum value would be chosen, which represents
 192 the closest point to the camera. If no measurement of Z is available for (x, y) , this would be treated as a missing
 193 value. To produce a localised search range $[d_{min}, d_{max}]$ for stereo matching, we used a patch centring around
 194 every pixel in the colour image to compute the minimum and maximum depth values. Denote (x, y, Z) as $Z_{(x,y)}$
 195 and the patch as $Z_{(\mathbf{m},\mathbf{n})}$,

$$|\mathbf{m} - x| \leq r, \quad |\mathbf{n} - y| \leq r \quad (13)$$

196 In effect, this allows mis-alignment up to r pixels when transforming the ToF image to the colour image. The
 197 maximum and minimum depths are then converted into disparities as,

$$d_{min(x,y)} = s f / \max\{Z_{(\mathbf{m},\mathbf{n})}\} - k \quad (14)$$

$$d_{max(x,y)} = s f / \min\{Z_{(\mathbf{m},\mathbf{n})}\} + k \quad (15)$$

198 The search range is expanded by k pixels (normally $0 \leq k \leq 3$) at each direction to allow for the noise in the
 199 ToF estimates. Given a localised search range $[d_{min}, d_{max}]$ for every pixel, a stereo method can then be used
 200 to find correspondences between images. In this work, for the data term cost D in (10), if $d_{(x,y)}$ is outside

201 the search range $[d_{min}, d_{max}]$ or $d_{(x,y)}$ is linked to a pixel outside the image, $D(d_{(x,y)})$ is set to the maximum
 202 pixel difference value T_d . If the localised search range $[d_{min}, d_{max}]$ is missing, $D(d_{(x,y)})$ is computed for all the
 203 possible disparities same as a dense stereo method.

204 It should also be noted that the search range is small at the object centre due to small depth variation and ToF
 205 measurements contribute more to the results, while the opposite can be observed at the depth discontinuities.

206 4.3 Quality score

207 Since pixel-by-pixel depth data are unavailable as ground truth for our images and many other applications,
 208 we propose a quantitative method accounting for the surface smoothness and the edge sharpness to evaluate
 209 estimation results. Leaf has depth edges along its boundary as seen in Fig. 5 and we can label depth edges to
 210 quantify how well the result has preserved them. Leaf boundaries shown in Fig. 5 were obtained manually. We
 211 only performed this manual edge labelling at this evaluation stage to produce ground truth for depth edges,
 212 and neither the ToF transformation nor the stereo method required any intervention after calibration. The area
 213 within the leaf boundaries is considered a leaf surface, and the final output consists of two binary images, a
 214 surface and an edge image. We applied the Canny edge filter [13] with non-maximum suppression to compute
 215 the smoothness of the surface and sharpness of the depth edges. Surface smoothness penalty P_s was calculated
 216 for the surface image as,

$$P_s = \overline{M(d)_{(x_s, y_s)}} \quad (16)$$

217 where M represents the Canny edge filter using a Gaussian that has a standard deviation of 1 and a radius of
 218 1.5 for non-maximum suppression. (x_s, y_s) are surface pixels. Edge sharpness score S_e was calculated for the
 219 edge image as,

$$S_e = \overline{(g * M(d))_{(x_e, y_e)}} \quad (17)$$

220 where $*$ is the 2-dimensional convolution operation [13] and g denotes a Gaussian filter in order to deal with
 221 thin and sharp depth edges. In this work, we set the neighbourhood size of the Gaussian filter to 15 and the

Table 1 Steps for automatically detecting frontal leaves.

1. Select the disparity plane nearest to the camera $\max\{d\}$.
2. Select a point (x_c, y_c) from $\max\{d\}$ that has the minimum value of combined edge magnitude $M(I, d, \gamma)$.
3. Use (x_c, y_c) as the seed point, perform region growing method to segment a leaf (x_l, y_l) .
4. Set $d_{(x_l, y_l)}$ to the furthest to the camera $\min\{d\}$.
5. Repeat 1 - 4 for the next leaf.

222 standard deviation to 5. The effects of (16) and (17) can be seen in Fig. 5. A quality score S accounting for the
 223 surface smoothness P_s and the edge sharpness S_e was computed as below,

$$S = S_e - P_s \quad (18)$$

224 The score S penalises displacement between defined depth edges and depth edges by a dense method while
 225 requiring the surface to be smooth. S is a relative score that becomes meaningful when comparing two dense
 226 methods. Given S calculated for two methods, the stereo result with the lower S score has more blurry depth
 227 edges (smaller S_e), more noise on surface (larger P_s), or both. Consequently, for our application in this paper,
 228 a dense method with a higher S score is preferred over one with a lower S score.

229 It should be noted that Sobel edge magnitude can also be used for M to calculate P_s and S_e as in [42].
 230 However, automatic leaf detection described in section 5 uses the same Canny edge filter.

231 4.4 Parameter tuning

232 One way of using the quality score in (18) is tuning parameters. Let θ be a set of parameters that is required
 233 by a method (e.g. our method combining stereo and ToF or other stereo methods), and the sum of score $\sum S$
 234 on some objects with manually labelled depth edges is used as a quality measure. In this work, we manually
 235 labelled three leaves in calibration images. The set of parameters θ that maximises $\sum S$ is considered as the
 236 tuned parameter set. Since θ has one or more parameters, we sequentially optimise all possible values in each
 237 parameter leading to a local optimisation $\hat{\theta}$. Several iterations defined by the number of parameters in θ are
 238 then used to refine $\hat{\theta}$, and the order of parameters is randomised in each iteration.

239 5 Leaf Detection

240 Given a colour image I and its corresponding disparity image d , this section describes a general approach for
 241 automatically extracting leaf boundaries from images. Let M represent a filter that produces edge magnitude,
 242 and $M(I)$ and $M(d)$ are the edge magnitude on the colour image I and the disparity image d respectively. It is
 243 possible to use either $M(I)$ or $M(d)$ alone for boundary detection. However, both $M(I)$ and $M(d)$ are compli-
 244 cated and noisy as shown in Fig. 6, which is hard to produce reasonable boundary detection. A combination of
 245 $M(I)$ and $M(d)$ can simplify the problem, which is obtained as follows,

$$M(I, d, \gamma) = \gamma M(I) + (1 - \gamma) M(d) \quad (19)$$

246 where γ is a weighting coefficient and $0 \leq \gamma \leq 1$. The idea is based on the fact that object boundaries could be
 247 enhanced by blending edges existing in the colour and disparity images, despite one of them could be weak in
 248 edge magnitude.

249 We assumed that some leaves are in foreground closer to the camera than other objects like stems, since
 250 leaves have to reach out for maximising light interception. Based on this assumption and the combined edge
 251 magnitude $M(I, d, \gamma)$, an automated procedure was then developed to detect the nearest frontal leaf. A summary
 252 of steps in the procedure is shown in Table 1. The second step attempts to select a point on a leaf surface as the
 253 seed point, which should have a small value of combined edge magnitude $M(I, d, \gamma)$. We used a region growing
 254 method [15] on edge magnitude produced by the filter M to extract regions, since the method is well-understood
 255 and performs well with respect to noise. The criterion we used was to compare edge magnitude of the adjacent
 256 pixels near the region borders to the region's mean value. A thresholding parameter T_c acting as a similarity
 257 threshold value was used to determine the terminal condition. There are also a number of other alternatives for
 258 foreground extraction using ToF and colour cameras [45, 8, 16, 40].

259 In this paper, we have found that the Canny edge filter with non-maximum suppression described in section
 260 4.3 can be used for the filter M . The Canny edge filter was configured to use a Gaussian that has a standard
 261 deviation of 1 and a radius of 1.5 for non-maximum suppression. For calculating edge magnitude on colour
 262 image $M(I)$, I is transformed from RGB to CIELAB colour space, and the edge magnitude by the Canny filter

263 is separately calculated for each component and is subsequently summed up. There are other choices of M that
264 can be used (e.g. [3]), and more methods can be found in the Berkeley segmentation benchmark [33].

265 Fig. 6(e) shows the output of (19) using the Canny filter, and $M(I, d, \gamma)$ emphasises edges found in both
266 $M(I)$ and $M(d)$ that corresponds to the boundary of our interest. Fig. 7 shows automatic leaf detection using
267 our method.

268 6 Surface Reconstruction

269 Given an identified leaf (x_l, y_l) in an image with colour and disparity data d , it is possible to reconstruct and
270 measure the leaf shape in 3D by a triangular mesh representation. Vertices in the mesh correspond to pixels
271 in the image which is $(x_l, y_l, d_{(x_l, y_l)})$, and the edges in the mesh are built by connecting nearest neighbouring
272 vertices. The use of a triangular mesh representation allows calculating the surface area in 3D analogous to
273 manual measurement. $(x_l, y_l, d_{(x_l, y_l)})$ is transformed to the world coordinate (X_l, Y_l, Z_l) as shown in section 3,
274 and areas of all triangles are then summed up.

275 An immediate problem of the triangular mesh representation is the ‘rice terrace’ effect as shown in Fig. 8(a).
276 This aliasing effect was caused by discretisation in the depth estimates at such a small scale. In this work,
277 baseline s is small for stereo matching and disparity estimates are integers, which cannot cope with the demand
278 for an accurate reconstruction. In addition to the inaccurate visual reconstruction, the ‘rice terrace’ effect would
279 over-estimate the surface area that is not desirable.

280 One approach is to develop sub-pixel accuracy stereo as described by [39], but in this work we propose a
281 method to smooth the depth data Z_l at the reconstruction stage after extracting leaves by combining stereo
282 and ToF. We decided to use local regression (LOESS) [6, 30], and it is based on the idea that any function can
283 be well approximated in a small neighbourhood by a low-order polynomial and that simple models can easily be
284 fit to data. A linear LOESS model was used in this work, and it requires a specific smoothing parameter β that
285 is a percentage of the total number of data points. The effect of LOESS smoothing increases with increasing β
286 as shown in Fig. 8. Fig. 8(d) and 8(e) also show the histogram of the residual image to illustrate the effect of
287 smoothing.

288 β can be defined empirically or by Generalised Cross Validation (GCV) [14]. However, we prefer a smooth
289 surface which would have larger residual on the flat planes of the ‘rice terraces’ rather than no or small residual,

290 and GCV is not built for correlated errors. An ad-hoc procedure to determine the smoothing parameter β
 291 is therefore developed. Since we would like to eliminate the aliasing effects, the residuals would not have a
 292 pronounced peak at 0 in histogram as in Fig. 8(d). Denote the histogram count as h_c with bin width b . The
 293 number of residuals at 0 is $h_c(0)$, and its two adjacent bins of the histogram are $h_c(-b)$ and $h_c(b)$ respectively.
 294 The stopping criterion for increasing the smoothing parameter is,

$$\min \left\{ h_{(0)} - h_{(-b)}, h_{(0)} - h_{(b)} \right\} \frac{1}{h_{(0)}} < T_h \quad (20)$$

295 where T_h is a parameter in percentage that controls the difference of histogram counts between the residual at
 296 0 and its two adjacent bins.

297 7 Results

298 7.1 Depth estimation

299 We compare three dense stereo algorithms with our method on some challenging pepper plant images. Images
 300 used for parameter tuning are calibration images, and the others are used as validation images. Once parameter
 301 tuning has been performed on calibration images as in section 4.4, none of the methods require intervention in
 302 the validation step.

303 We present results on three sets of images capturing well-developed plants (Plant 1 - 3). Plant 1 acts as the
 304 calibration image, and Plant 2 and 3 are validation images. Let SIFTflow, Shape and GC represent methods by
 305 Liu *et. al.* [29], Ogale and Aloimonos [34] and Boykov *et. al.* [5] respectively. GC refers to the graph cut method
 306 without using ToF, and GC+ToF is our method. By parameter tuning described in section 4.4, SIFTflow was
 307 configured with a 5-level pyramid, 5×5 window, $\alpha = 1$ and $\gamma = 0.001$. The α in the Shape method was set
 308 to 2. Parameters $T_e, T_d, T_k, \alpha_v, n_v$ for GC were set as 25, 20, 6, 4, 4. For GC+ToF, the same parameters for GC
 309 were used for dense stereo and ToF parameters r and k were set to 10 and 1 respectively. Since these methods
 310 are established, readers can see the effects of these parameters by following [29,34,5] for SIFTflow, Shape and
 311 GC respectively.

312 Fig. 9 shows qualitative stereo results produced by the four methods on Plant 1. Methods GC and GC+ToF
 313 produced results with leaves recognisable from the background. SIFTflow produced smooth results but did not

Table 2 Numerical summary of quality evaluation for Leaf 1, Leaf 2 and Leaf 3. S_e refers to edge sharpness, P_s refers to surface smoothness and S is the quality score. For our application in this paper, we would like to find a method giving the highest S score.

	Leaf 1			Leaf 2			Leaf 3		
	S_e	P_s	S	S_e	P_s	S	S_e	P_s	S
SIFTflow	0.22	0.10	0.12	0.11	0.01	0.10	0.04	0.02	0.02
Shape	0.70	0.44	0.26	0.74	0.20	0.54	0.27	0.10	0.18
GC	0.49	0.20	0.29	0.87	0.04	0.83	0.20	0.04	0.16
GC+ToF	1.34	0.18	1.15	1.31	0.05	1.26	0.36	0.06	0.30

Table 3 Quantitative summary of quality scores for both calibration and validation images. Figures shown here are total quality scores for three leaves in a image, $\sum S$.

	Calibration	Validation	
	Plant 1	Plant 2	Plant 3
SIFTflow	0.24	0.18	0.36
Shape	0.98	1.32	0.45
GC	1.28	1.29	0.51
GC+ToF	2.71	2.04	0.63

314 preserve discontinuity, while Shape showed the opposite effect. This can be further examined in Fig. 10, which
 315 shows a closer view of three leaves. The edge was weak for SIFTflow, although the surface was the most smooth.
 316 Method Shape suffered from noises on the surface, and GC failed to produce some depth edges. In comparison,
 317 GC+ToF produced the best qualitative results among the four methods.

318 A summary of quantitative results (S_e, P_s, S) for all three leaves is shown in Table 2. Similar to the findings
 319 in the qualitative results above, we see that GC+ToF produced sharp depth edges represented by a high S_e
 320 score especially for Leaf 1 and Leaf 2. The ranking of methods produced by the score S is also consistent with
 321 the qualitative results for the two leaves. Leaf 3 is in front of another leaf, and the magnitude of depth edges is
 322 therefore not as strong as those in Leaf 1 and Leaf 2. GC+ToF produced the best scores S_e and S among the
 323 four methods.

324 This section has shown results on the calibration image (Plant 1) to illustrate the behaviour of the four
 325 methods. Furthermore, two sets of results on validation images (Plant 2 - 3) have been produced. Table 3
 326 presents $\sum S$ for Plant 1 - 3. By using ToF as a localised search range, the estimation results were improved by
 327 at least 23% measured by the score $\sum S$.

328 It should be noted that the leaf boundary was manually selected here for comparison purposes, since auto-
 329 matic leaf detection can be difficult for estimates in Fig. 9. The next section will present the results of automatic
 330 leaf detection using our method.

331 7.2 Leaf detection and area measurements

332 A validation trial using 44 experimental plants (11 plots of 11 genotypes, each plot also has four border plants)
333 was carried out to provide a set of validation images. The 44 experimental plants grew in a standard double-
334 row arrangement, with 22 of them visible from each side of the row by our setup. The validation images were
335 collected first, and four leaves in the foreground of images were then manually detached from every genotype
336 giving 88 leaves in total. The positions of these leaves were annotated on the validation images in order to relate
337 to manual measurements as the ground truth. There were over 600 colour images to annotate, and these 88
338 leaves were identified in 244 images as 248 separate measurements. We ignored these leaves spanning across two
339 images (i.e. partial view of a complete leaf). Finally, the manual measurements of these 88 leaves were obtained
340 by removing each leaf from its plant and scanning using an industry-precision LI-COR 3100 leaf area meter.

341 Calibration images for parameter tuning were collected in a different trial. The parameters determined in
342 section 7.1 were used for our GC+ToF method. Based on qualitative results of calibration images, parameters T_c
343 and γ in leaf detection were set to 0.5 and 0.4, and three leaves were automatically detected in every calibration
344 image. Parameters b and T_h in surface reconstruction were empirically set to 0.2 and 0.1.

345 Our methods successfully extracted three leaves in each image producing 732 separate measurements from
346 244 images, but only 149 separate measurements on 59 leaves could be linked with the ground truth. This was
347 because the ground truth data were created first, and our methods were fully automated for analysing validation
348 images without accessing any manual annotation. Fig. 11 presents 15 examples randomly selected from the 149
349 boundary results.

350 Fig. 12 and Table 4 show the validation results of the 149 automatic measurements against 59 manual
351 measurements with average area of 102.0cm^2 . If no smoothing was applied, a lower correlation score and a
352 considerable larger RMSE value were obtained, which is due to the ‘rice terrace’ effect as expected. Using our
353 proposed smoothing method, the correlation between automatic and manual measurements is 0.97 and the
354 RMSE value is 10.97 cm^2 .

355 By averaging the 149 measurements for 59 leaves from different views, the correlation score has increased
356 to 0.98, and the RMSE value is reduced to 9.50 cm^2 , i.e. 9.3% of the average leaf area. These estimates have

Table 4 Correlation and Root-Mean-Square-Error (RMSE) results of 149 automatic leaf area measurements in validation data. Pearson correlation coefficient is used in this paper.

	Correlation	RMSE (cm ²)
No Smoothing	0.72	129.45
With Smoothing	0.97	10.97

357 been found to be of sufficient accuracy for plant breeding by identifying QTLs: positions on chromosomes which
 358 correlate significantly with measurements [20].

359 8 Discussion

360 Current practice is to measure leaf areas manually by destroying plants, which is also very costly in human time.
 361 Our method is non-destructive, and took about 3 minutes to record all images in a single row of plants. Then,
 362 average CPU times were 61sec/image for depth estimation, 1sec/image for leaf detection and 44sec/image for
 363 leaf area estimation. However, as images were processed off-line this is not critical, and could be speeded up by
 364 using more sophisticated or approximate algorithms or parallel processing.

365 The proposed 3D approach allows automatic measurement of the sizes of pepper leaves in a greenhouse.
 366 The setup could be extended to measure leaf sizes of other greenhouse crops such as cucumber or tomato, and
 367 in other situations, such as pot plants on a conveyer belt system. Similar approaches could be developed to
 368 measure size, orientation and shape of other components of plants, e.g. flowers, fruits, stems, and to exploit
 369 other multisensor systems [10,12,36].

370 Using stereo vision alone to extract individual leaves in scenes of plant structures with many leaves is clearly
 371 not sufficient. For occlusions and areas affected by unpredictable illumination, the data term in a global stereo
 372 framework (i.e. D in (10)) produces inaccurate energy costs since corresponding pixels are either unavailable
 373 or difficult to match. The minimised energy therefore does not represent the desirable results. Using ToF in
 374 these situations provides an estimate and reduces ambiguities. Another advantage is that dense stereo can be
 375 a super resolution technique for ToF images as discussed by [11,22,41]. Even using the latest ToF camera, the
 376 resolution of a ToF camera (320×240) is low compared with a colour camera, and this leads to errors at depth
 377 discontinuities. On the other hand, stereo vision is known for underperforming in low-textured areas usually
 378 at the object centres, but it can preserve discontinuity as we have presented in this paper (e.g. Fig. 10). Since

379 stereo and ToF complement each other, our methods to combine them can therefore be used for applications
380 aiming to improve the accuracy of measurement.

381 For our application, edge-preserving disparity result is very important. Since there are many leaves in an
382 image (Fig. 9) and they have the similar appearance as well as large intra-class variation in visual images, it is
383 impossible to extract one leaf from the rest without combining disparity estimation into the colour image (see
384 Fig. 11). As shown in Fig. 11 and 12, our methods were fully automated for the validation images without any
385 manual input (e.g. annotation), and were therefore proven to be effective for extracting reasonable foreground
386 boundaries.

387 Our current assumption that leaves are in the foreground limits our methods to find only frontal leaves. We
388 successfully collected 732 separate leaf measurements for 11 different genotypes of plants, and [20] presented
389 further results using our methods for 151 genotypes. This was the first viable approach to collect a large number
390 of leaf measurements in a non-destructive manner. In our work, foreground objects near to the cameras were
391 all leaves. In a more general framework, a detection verification procedure, using classification methods such as
392 support vector machine [9], could be developed for detecting target objects that are not in foreground.

393 Fig. 12 and Table 4 demonstrate the consequence of the ‘rice terrace’ effect, and highlight the importance
394 of surface smoothing when calculating surface area. Although we presented our method for one pair of stereo
395 images and one ToF image, it is in principle rather straightforward to apply it to multiple colour images and
396 one ToF image, or even to multiple colour and ToF images. Both the correlation score and the RMSE value can
397 be improved by averaging over multiple views, and the use of two or more views should be adopted in practice
398 to reduce occlusions. We hope to build on the work in this paper for combining multiple colour and ToF images,
399 and Kim *et. al.* [24] have shown some promising results on this subject.

400 9 Conclusions

401 This paper has presented an automatic approach for non-destructively measuring leaf surface area. Agreement
402 with the ground truth of manually measured leaf areas, shown in Fig. 12, is good, and sufficient for plant breeding
403 purposes [20]. Unlike most existing methods requiring individual plants to be transported to a controlled imaging
404 environment, our work collects measurements from plants in their own growing environment. Three frontal

405 leaves in an image (twelve from four images) were automatically measured, and they produced an estimate of
406 the average leaf area that can be used in plant growth analysis[20].

407 We have demonstrated that combining stereo and ToF images leads to discontinuity-preserving results,
408 which enable algorithms like Canny filter and region growing segmentation to extract individual leaves from a
409 cluttered scene. We have also highlighted the importance of surface smoothing for calculating surface area, and
410 proposed an adaptive way to choose the smoothing parameter.

411 Our approach has produced promising results on 149 automatic leaf area measurements in the validation
412 data, which had a correlation score of 0.97 and a RMSE value of 10.97 cm² against the manual measurements.
413 By using multiple views for the 59 leaves, the RMSE value reduced to 9.50 cm², with a correlation of 0.98.

414 The idea of combining stereo and ToF images has been proven useful for 3D measurements, and our approach
415 could potentially be applied for combining other modalities of images with large difference in image resolutions
416 and camera positions. Moreover, our automated approach is a major step forward in relation to the current
417 destructive and laborious practice for measuring leaf area.

418 10 Acknowledgements

419 This work is part of the Smart tools for Prediction and Improvement of Crop Yield (SPICY) project supported
420 by the European Community and funded by the KBBE FP7 programme. (Grant agreement number KBBE-
421 2008-211347) We also acknowledge Scottish Government funding.

422 References

- 423 1. Alenya, G., Dellen, B., Torras, C.: 3d modelling of leaves from color and tof data for robotized plant measuring. In: International
424 Conference on Robotics and Automation, ICRA 2011, pp. 3408–3414 (2011)
- 425 2. Alimi, N., Bink, M., Dieleman, J., Nicola, M., Wubs, M., Heuvelink, E., Magan, J., Voorrips, R., Jansen, J., Rodrigues, P.,
426 Heijden, G., Vercauteren, A., Vuylsteke, M., Song, Y., Glasbey, C., Barocsi, A., Lefebvre, V., Palloix, A., Eeuwijk, F.: Genetic
427 and QTL analyses of yield and a set of physiological traits in pepper. *Euphytica* **190**, 181–201 (2013)
- 428 3. Arbelaez, P., Maire, M., Fowlkes, C., Malik, J.: Contour detection and hierarchical image segmentation. *IEEE Trans. Pattern*
429 *Anal. Mach. Intell.* **33**, 898–916 (2011)
- 430 4. Beder, C., Bartczak, B., Koch, R.: A combined approach for estimating patchlets from pmd depth images and stereo intensity
431 images. In: Proceedings of the 29th DAGM conference on Pattern recognition, pp. 11–20. Springer-Verlag (2007)

- 432 5. Boykov, Y., Veksler, O., Zabih, R.: Fast approximate energy minimization via graph cuts. *IEEE Trans. Pattern Anal. Mach.*
433 *Intell.* **23**, 1222–1239 (2001)
- 434 6. Cleveland, W.S., Loader, C.L.: Smoothing by local regression: Principles and methods. In: W. Haerdle, M.G. Schimek (eds.)
435 *Statistical Theory and Computational Aspects of Smoothing*, pp. 10–49. Springer, New York (1996)
- 436 7. Collignon, A., Maes, F., Delaere, D., Vandermeulen, D., Suetens, P., Marchal, G.: Automated multi-modality image registration
437 based on information theory. *Information Processing in Medical Imaging* pp. 263–274 (1995)
- 438 8. Crabb, R., Tracey, C., Puranik, A., Davis, J.: Real-time foreground segmentation via range and color imaging. In: *IEEE*
439 *Computer Society Conference on Computer Vision and Pattern Recognition Workshops*, pp. 1–5 (2008)
- 440 9. Cristianini, N., Shawe-Taylor, J.: *An Introduction to Support Vector Machines and Other Kernel-based Learning Methods*.
441 Cambridge University Press (1999)
- 442 10. Dhondt, S., Wuyts, N., Inze, D.: Cell to whole-plant phenotyping: the best is yet to come. *Trends in Plant Science* (in press)
- 443 11. Diebel, J., Thrun, S.: An application of markov random fields to range sensing. In: *Proceedings of Conference on Neural*
444 *Information Processing Systems (NIPS)*, pp. 291–298. MIT Press, Cambridge, MA (2005)
- 445 12. Fiorani, F., Schurr, U.: Future scenarios for plant phenotyping. *Plant Biology* **64**, 267–291 (2013)
- 446 13. Glasbey, C.A., Horgan, G.W.: *Image analysis for the biological sciences*. John Wiley & Sons, Inc., New York, NY, USA (1995)
- 447 14. Golub, G.H., Heath, M., Wahba, G.: Generalized cross-validation as a method for choosing a good ridge parameter. *Techno-*
448 *metrics* **21**(2), pp. 215–223 (1979)
- 449 15. Gonzalez, R.C., Woods, R.E.: *Digital Image Processing*, 3rd edn. Prentice Hall (2008)
- 450 16. Gudmundsson, S.A., , Pardis, M., Casas, J.R., Sveinsson, J.R., Aanaes, H., Larsen, R.: Improved 3d reconstruction in smart-
451 room environments using tof imaging. *Comput. Vis. Image Underst.* **114**, 1376–1384 (2010)
- 452 17. Gudmundsson, S.A., Aanaes, H., Larsen, R.: Fusion of stereo vision and time-of-flight imaging for improved 3d estimation.
453 *International Journal of Intelligent Systems Technologies and Applications* **5**(3/4), 425–433 (2008)
- 454 18. Hahne, U., Alexa, M.: Combining time-of-flight depth and stereo images without accurate extrinsic calibration. *International*
455 *Journal of Intelligent Systems Technologies and Applications* **5**(3/4), 325–333 (2008)
- 456 19. Hartley, R.I., Zisserman, A.: *Multiple View Geometry in Computer Vision*, second edn. Cambridge University Press (2004)
- 457 20. van der Heijden, G., Song, Y., Horgan, G., Polder, G., Dieleman, A., Bink, M., Palloix, A., van Eeuwijk, F., Glasbey, C.:
458 SPICY: towards automated phenotyping of large pepper plants in the greenhouse. *Functional Plant Biology* **39**(11), 870–877
459 (2012)
- 460 21. Heuvelink, E.: Evaluation of a dynamic simulation model for tomato crop growth and development. *Annals of Botany* **83**,
461 413–422 (1998)
- 462 22. Huhle, B., Schairer, T., Jenke, P., Straier, W.: Fusion of range and color images for denoising and resolution enhancement with
463 a non-local filter. *Comput. Vis. Image Underst.* **114**, 1336–1345 (2010)
- 464 23. Jonckheere, I., Fleck, S., Nackaerts, K., Muys, B., Coppin, P., Weiss, M., Baret, F.: Review of methods for in situ leaf area
465 index determination: Part i. theories, sensors and hemispherical photography. *Agricultural and Forest Meteorology* **121**(1-2),
466 19 – 35 (2004)

- 467 24. Kim, Y., Theobalt, C., Diebel, J., Kosecka, J., Micusik, B., Thrun, S.: Multi-view image and tof sensor fusion for dense 3d
468 reconstruction. In: Proceedings of IEEE 12th International Conference on Computer Vision Workshops, ICCV Workshops, pp.
469 1542–1549 (2009)
- 470 25. Kolb, A., Barth, E., Koch, R., Larsen, R.: Time-of-Flight Sensors in Computer Graphics. In: M. Pauly, G. Greiner (eds.)
471 Eurographics 2009 - State of the Art Reports, pp. 119–134. Eurographics (2009)
- 472 26. Kumar, N., Belhumeur, P.N., Biswas, A., Jacobs, D.W., Kress, W.J., Lopez, I.C., Soares, J.V.B.: Leafsnap: a computer vision
473 system for automatic plant species identification. In: Computer Vision – ECCV 2012, pp 502–516. (2012)
- 474 27. Lindner, M., Lambers, M., Kolb, A.: Sub-pixel data fusion and edge-enhanced distance refinement for 2d/3d images. Interna-
475 tional Journal of Intelligent Systems Technologies and Applications **5**, 344–354 (2008)
- 476 28. Lindner, M., Schiller, I., Kolb, A., Koch, R.: Time-of-flight sensor calibration for accurate range sensing. *Comput. Vis. Image*
477 *Underst.* **114**, 1318–1328 (2010)
- 478 29. Liu, C., Yuen, J., Torralba, A.: Sift flow: Dense correspondence across scenes and its applications. *IEEE Trans. Pattern Anal.*
479 *Mach. Intell.* (2010)
- 480 30. Loader, C.: Local regression and likelihood (2001).
481 <http://cm.bell-labs.com/stat/project/locfit>
- 482 31. Lowe, D.: Object recognition from local scale-invariant features. In: Proceedings of the Seventh IEEE International Conference
483 on Computer Vision, vol. 2, pp. 1150–1157 (1999)
- 484 32. Marcelis, L., Heuvelink, E., Goudriaan, J.: Modelling biomass production and yield of horticultural crops: a review. *Scientia*
485 *Horticulturae* **74**(1-2), 83 – 111 (1998)
- 486 33. Martin, D., Fowlkes, C., Tal, D., Malik, J.: A database of human segmented natural images and its application to evaluating
487 segmentation algorithms and measuring ecological statistics. In: Proc. 8th Int’l Conf. Computer Vision, vol. 2, pp. 416–423
488 (2001)
- 489 34. Ogale, A.S., Aloimonos, Y.: Shape and the stereo correspondence problem. *Int. J. Comput. Vision* **65**, 147–162 (2005)
- 490 35. Parsons, N., Edmondson, R., Song, Y.: Image analysis and statistical modelling for measurement and quality assessment of
491 ornamental horticulture crops in glasshouses. *Biosystems Engineering* **104**(2), 161 – 168 (2009)
- 492 36. Pieruschka, R., Poorter, H.: Phenotyping plants: genes, phenes and machines. *Functional Plant Biology* **39**, 813–820 (2012)
- 493 37. Polder, G., van der Heijden, G.W.A.M., Glasbey, C.A., Song, Y., Dieleman, J.A.: Spy-See - Advanced vision system for
494 phenotyping in greenhouses. In: Proceedings of the MINET Conference: Measurement, sensation and cognition, pp. 115–117.
495 National Physical Laboratory (2009)
- 496 38. Rajendran, K., Tester, M., Roy, S.J.: Quantifying the three main components of salinity tolerance in cereals. *Plant, Cell and*
497 *Environment* **32**(3), 237–249 (2009)
- 498 39. Scharstein, D., Szeliski, R.: A taxonomy and evaluation of dense two-frame stereo correspondence algorithms. *Int. J. Comput.*
499 *Vision* **47**(1-3), 7–42 (2002)
- 500 40. Schiller, I., Koch, R.: Improved video segmentation by adaptive combination of depth keying and mixture-of-gaussians. In:
501 The 17th Scandinavian Conference on Image Analysis, SCIA 2011, pp. 59–68. Ystad, Sweden (2011)

- 502 41. Schuon, S., Theobalt, C., Davis, J., Thrun, S.: Lidarboost: Depth superresolution for tof 3d shape scanning. In: Proceedings
503 of the IEEE CVPR 2009, pp. 343 – 350 (2009)
- 504 42. Song, Y., Glasbey, C.A., van der Heijden, G.W.A.M., Polder, G., Dieleman, J.A.: Combining stereo and Time-of-Flight images
505 with application to automatic plant phenotyping. In: The 17th Scandinavian Conference on Image Analysis, SCIA 2011, pp.
506 467–478. Ystad, Sweden (2011)
- 507 43. Song, Y., Wilson, R., Edmondson, R., Parsons, N.: Surface modelling of plants from stereo images. In: Proceedings of the 6th
508 International Conference on 3-D Digital Imaging and Modeling. Montreal, Canada (2007)
- 509 44. Tola, E., Lepetit, V., Fua, P.: Daisy: an efficient dense descriptor applied to wide baseline stereo. *IEEE Trans. Pattern Anal.*
510 *Mach. Intell.* **32**(5), 815–830 (2010)
- 511 45. Wang, L., Zhang, C., Yang, R., Zhang, C.: TofCut: Towards Robust Real-Time Foreground Extraction Using a Time-of-Flight
512 Camera. In: The Fifth International Symposium on 3D Data Processing, Visualization and Transmission, 3DPVT 2010 (2010)
- 513 46. Wubs, A.M., Ma, Y., Hemerik, L., Heuvelink, E.: Fruit set and yield patterns in six capsicum cultivars. *HortScience* **44**,
514 1296–1301 (2009)
- 515 47. Zhu, J., Wang, L., Yang, R., Davis, J.: Fusion of time-of-flight depth and stereo for high accuracy depth maps. In: Proceedings
516 of the IEEE CVPR 2008, pp. 1–8 (2008)

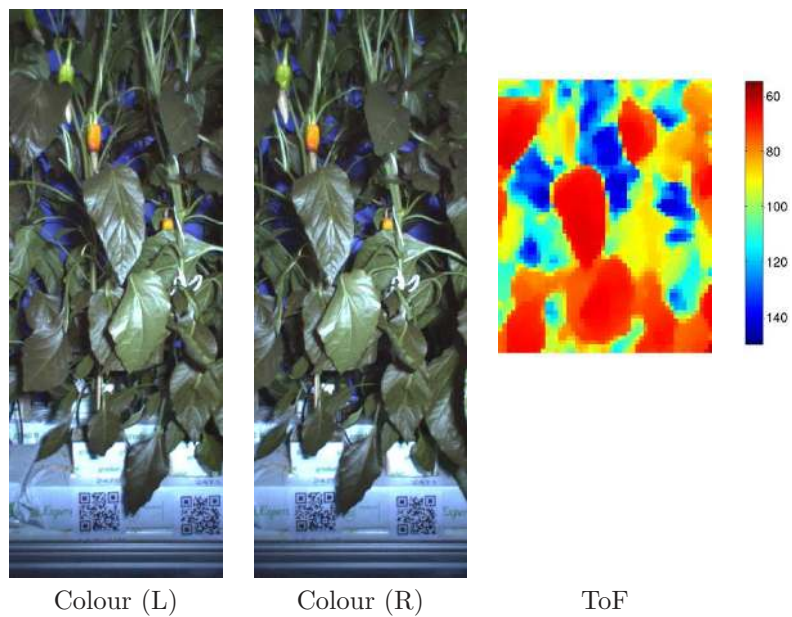


Fig. 1 Pepper plant images. Colour (L) and Colour (R) are a stereo pair of images of pepper plants, and Colour (R) is the base image. ToF is the depth image in cm matching the base image. The colour images are 480×1280 in size, while ToF image is only 64×48 . The ratio of pixels between colour and ToF images is 200 : 1.



(a)

(b)

Fig. 2 Our system has four camera rigs vertically stacked to capture pepper plants, and each one has colour and ToF cameras and a flash light as seen in (a). Our system also moves around in a greenhouse that is exposed to unpredictable lighting including the sun and reflection as in (b).

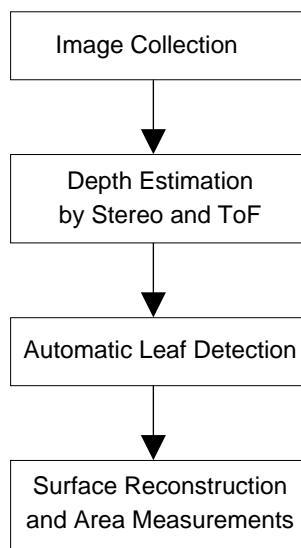


Fig. 3 Overview of main processes of our system for measuring leaf area.

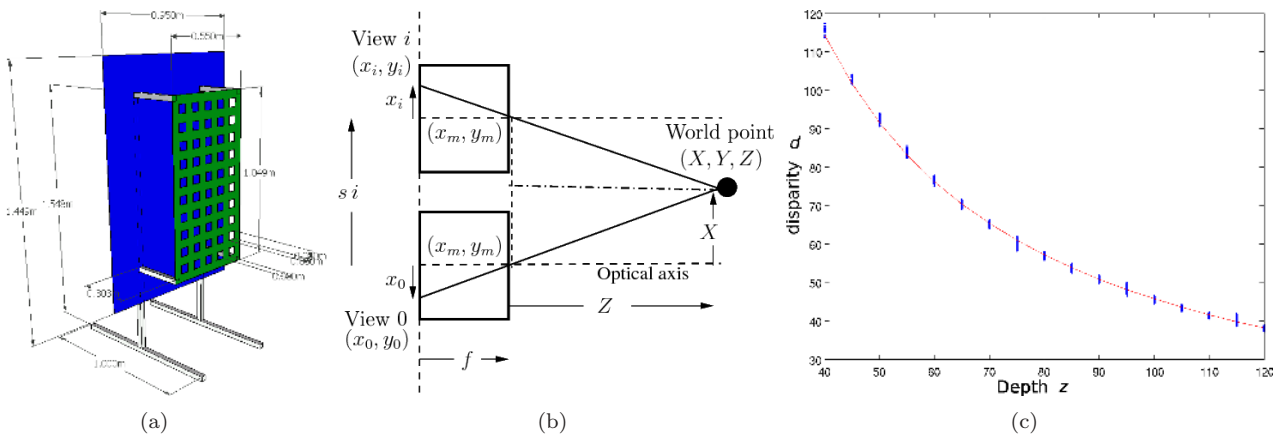


Fig. 4 Camera calibration: (a) diagram of calibration board; (b) diagram illustrating the transformation of a point between colour images and world coordinate. Dashed line represents the optical axis going through the principal point (x_m, y_m) . y_i and Y represent axis perpendicular to the page with the same projective properties. (c) plot of the relationship between depth Z in cm and disparities d in pixels for colour camera. Blue dots were disparity measurements d for each Z , and the red line was the fit by (6).

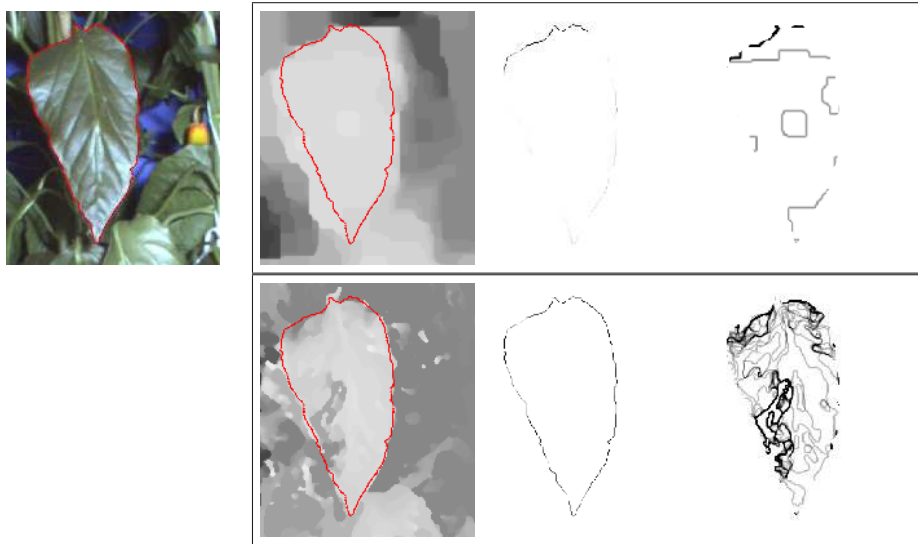


Fig. 5 Examples to illustrate quality scores on two disparity results. The colour image with depth edges plotted in red is shown on the left. The grey values in each panel represent the disparity (left), S_e (middle) and P_s (right) respectively. The disparity maps use a grey value scale of 20-90 pixels black-white. The S_e and P_s images also use a common scale.

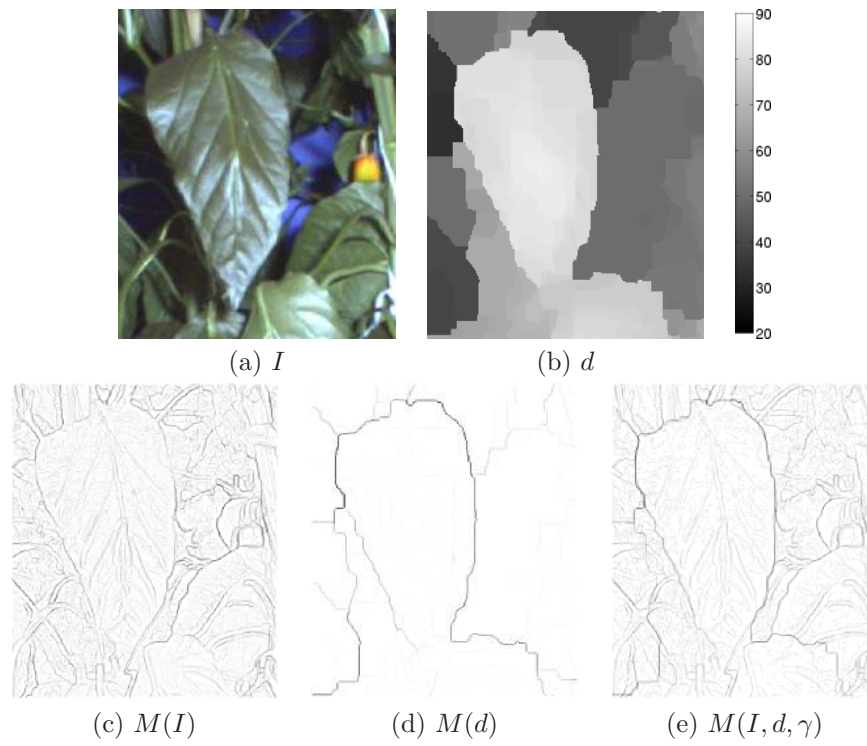


Fig. 6 Edge magnitude of colour and disparity images. (a) and (b) show the colour I and estimated disparity d images respectively. (c) and (d) show edge magnitudes by the Canny filter $M(I)$ and $M(d)$, and (e) shows the combined Canny edge magnitude. For this example, the weight coefficient γ is set to 0.4.



Fig. 7 Comparison between automatic leaf detection and manual selection. Yellow boundaries represent automatic leaf detection by our method and red boundaries represent manual selection.

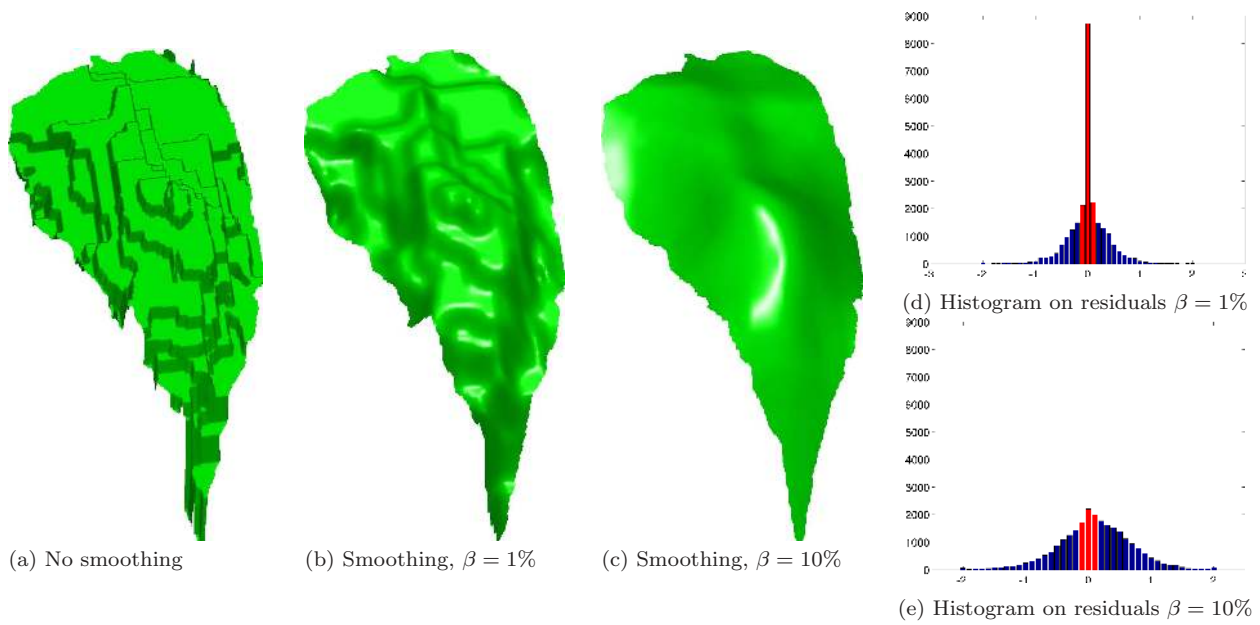


Fig. 8 Reconstruction of a leaf surface with smoothing parameter β . (a) shows reconstruction without smoothing and the ‘rice terrace’ effect is clear. (b) shows little LOESS smoothing with $\beta = 1\%$, which is also clear from histogram on residuals in (d). (c) shows LOESS smoothing with $\beta = 10\%$. (e) shows histogram on residuals with $\beta = 10\%$, and there is a pattern of changes in the three bins near 0 in red colour compared with (d).

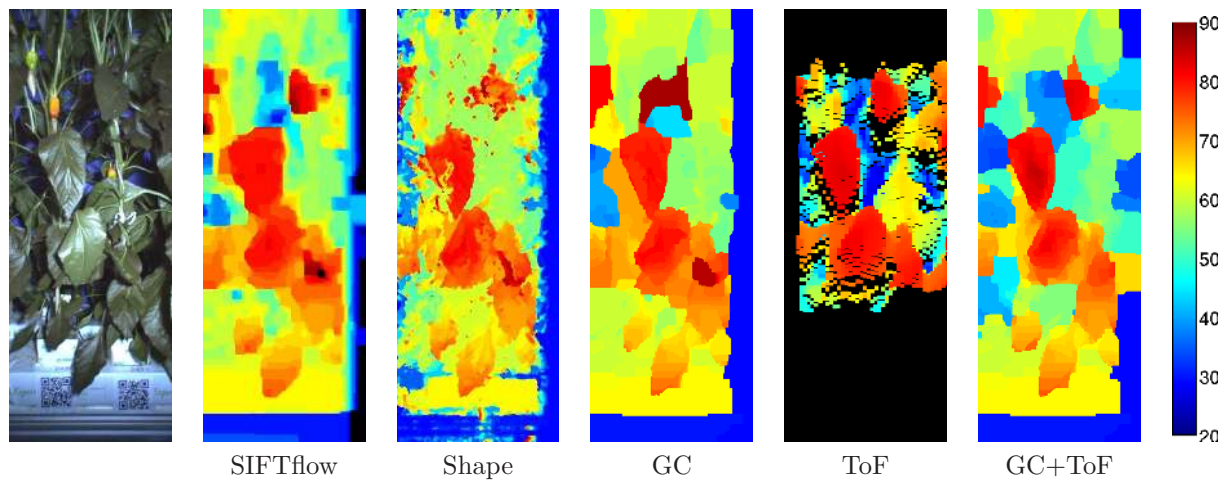


Fig. 9 Disparity results on the 'Plant 1'. SIFTflow, Shape and GC represent methods by Liu *et. al.* [29], Ogale and Aloimonos [34] and Boykov *et. al.* [5] respectively. ToF shows transformed points in colour image coordinates and the black pixels indicate missing ToF information. GC+ToF is our method combining stereo and ToF images.

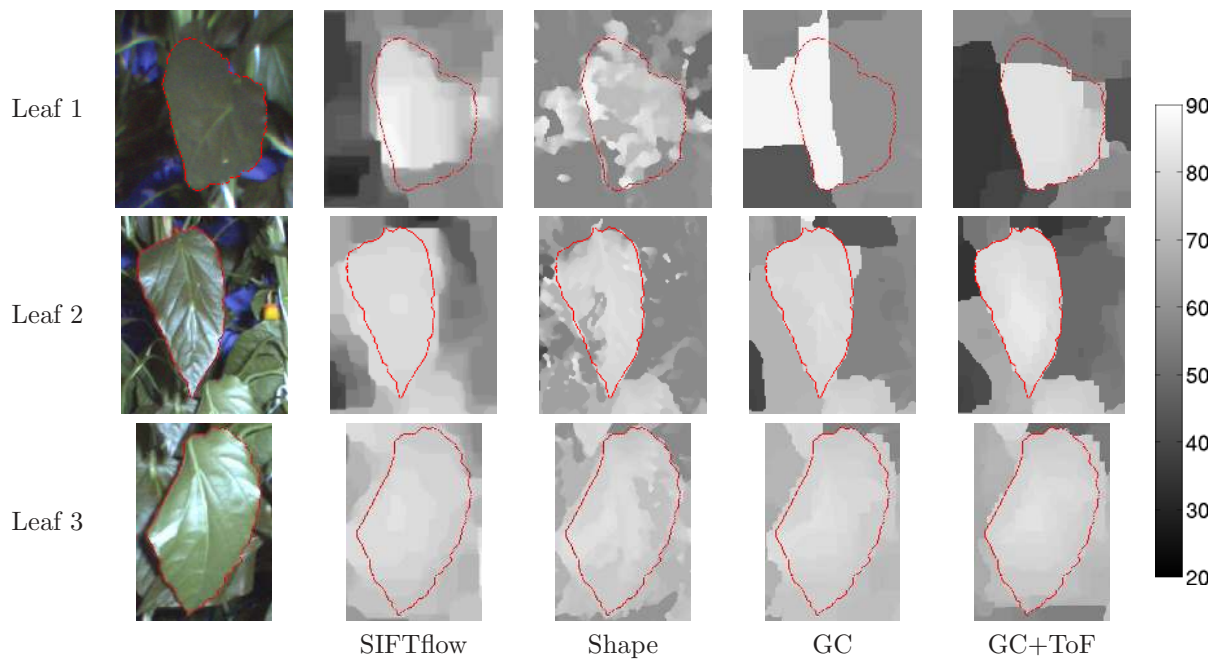


Fig. 10 Results for three leaves from top to bottom in Plant 1. Depth edges are plotted in red.

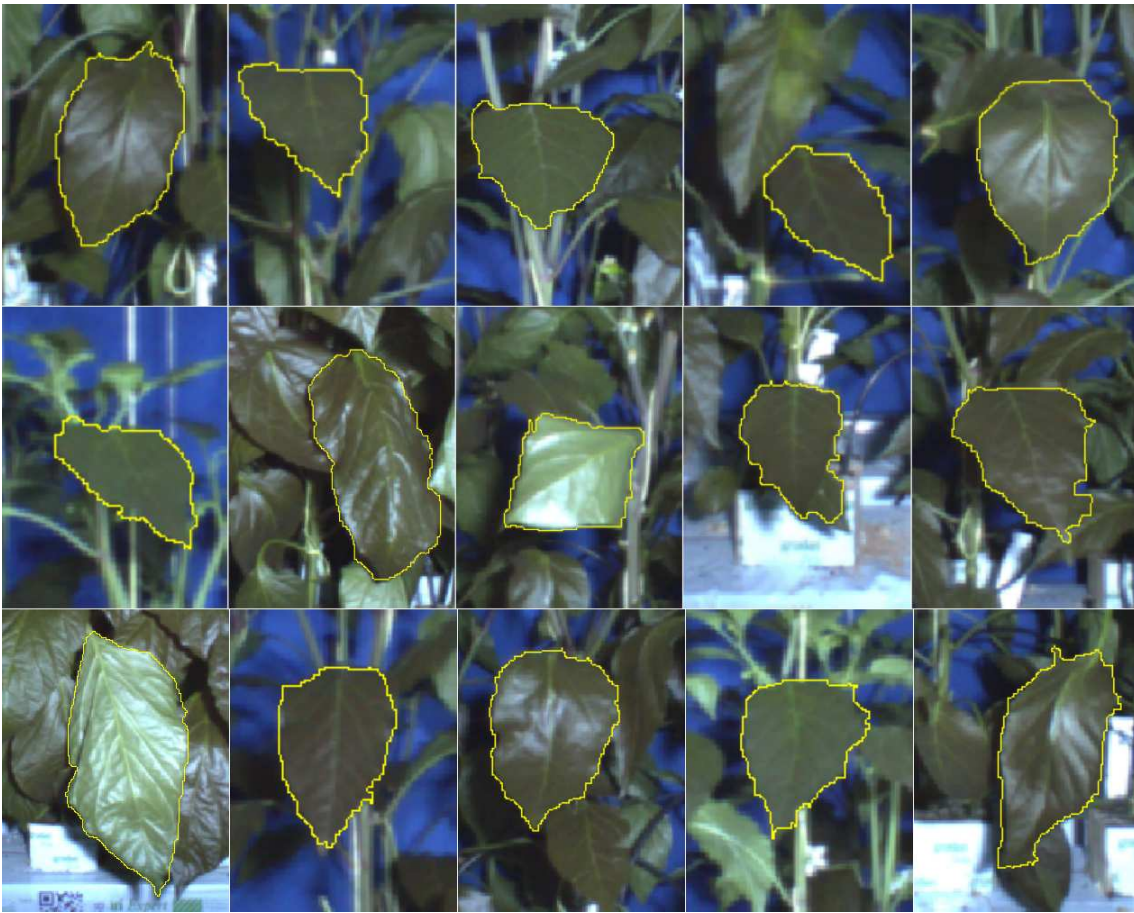


Fig. 11 Results on 15 examples of automatic leaf detection by our method. The yellow boundary outlines an automatically identified leaf. The 15 examples were randomly selected from the possible 149 leaves automatically detected in the validation images.

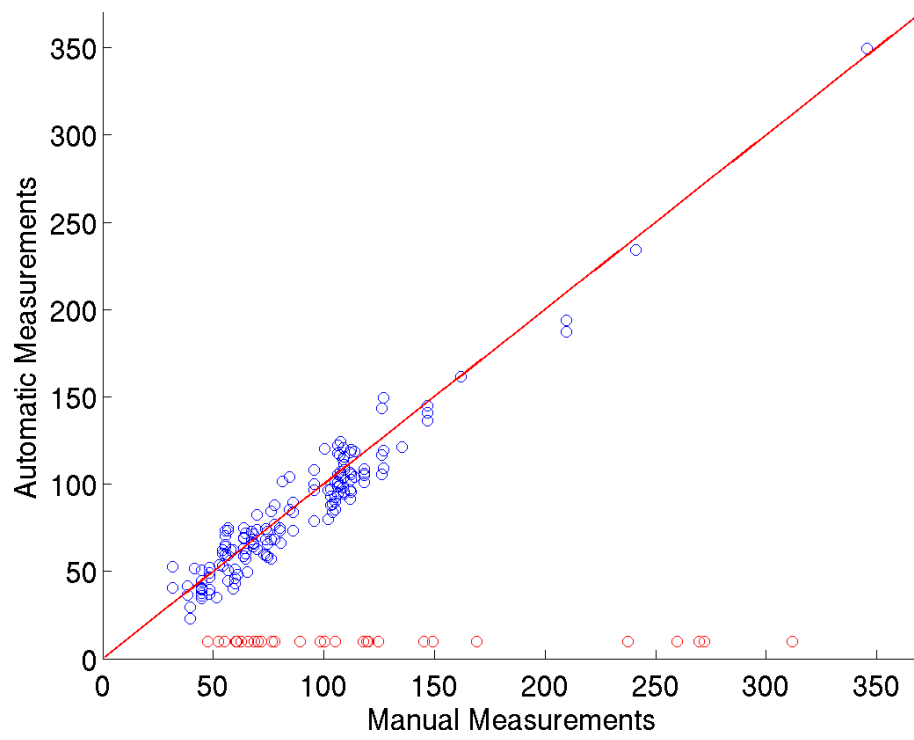


Fig. 12 Plots of 149 automatic leaf area measurements automatically obtained using our system against manual measurements on the validation images. The x-axis is manual measurements in cm², and the y-axis is automatic measurements in cm². The red line is the 1 : 1 reference line, and the red circles show the 29 leaves that could not be automatically detected.